

# Generalization Bounds for Domain Adaptation

Chao Zhang<sup>1</sup>, Lei Zhang<sup>2</sup>, Jieping Ye<sup>1,3</sup>

<sup>1</sup>Center for Evolutionary Medicine and Informatics, The Biodesign Institute,  
and <sup>3</sup>Computer Science and Engineering, Arizona State University, Tempe, USA

zhangchao1015@gmail.com; jieping.ye@asu.edu

<sup>2</sup>School of Computer Science and Technology,  
Nanjing University of Science and Technology, Nanjing, P.R. China  
zhanglei.njust@yahoo.com.cn

April 8, 2013

## Abstract

In this paper, we provide a new framework to obtain the generalization bounds of the learning process for domain adaptation, and then apply the derived bounds to analyze the asymptotical convergence of the learning process. Without loss of generality, we consider two kinds of representative domain adaptation: one is with multiple sources and the other is combining source and target data.

In particular, we use the integral probability metric to measure the difference between two domains. For either kind of domain adaptation, we develop a related Hoeffding-type deviation inequality and a symmetrization inequality to achieve the corresponding generalization bound based on the uniform entropy number. We also generalized the classical McDiarmid's inequality to a more general setting where independent random variables can take values from different domains. By using this inequality, we then obtain generalization bounds based on the Rademacher complexity. Afterwards, we analyze the asymptotic convergence and the rate of convergence of the learning process for such kind of domain adaptation. Meanwhile, we discuss the factors that affect the asymptotic behavior of the learning process and the numerical experiments support our theoretical findings as well.

## 1 Introduction

The generalization bound measures the probability that a function, chosen from a function class by an algorithm, has a sufficiently small error and plays an important role in statistical learning theory (see Vapnik, 1999; Bousquet *et al.*, 2004). The generalization bounds have been widely used to study the consistency of the ERM-based learning process (Vapnik, 1999), the asymptotic convergence of empirical process (Van der Vaart and Wellner, 1996) and the learnability of learning models (Blumer *et al.*, 1989). Generally, there are three essential aspects to obtain the generalization bounds of a specific learning process: complexity measures of function classes, deviation (or concentration) inequalities and symmetrization inequalities related to the learning process. For example, Van der Vaart and Wellner (1996) presented the generalization bounds

based on the Rademacher complexity and the covering number, respectively. Vapnik (1999) gave the generalization bounds based on the Vapnik-Chervonenkis (VC) dimension. Bartlett *et al.* (2005) proposed the local Rademacher complexity and obtained a sharp generalization bound for a particular function class  $\{f \in \mathcal{F} : \mathbb{E}f^2 < \beta \mathbb{E}f, \beta > 0\}$ . Hussain and Shawe-Taylor (2011) showed improved loss bounds for multiple kernel learning.

It is noteworthy that the aforementioned results of statistical learning theory are all built under the assumption that training and test data are drawn from the same distribution (or briefly called the assumption of same distribution). This assumption may not be valid in the situation that training and test data have different distributions, which will arise in many practical applications including speech recognition (Jiang and Zhai, 2007) and natural language processing (Blitzer *et al.*, 2007). Domain adaptation has recently been proposed to handle this situation and it is aimed to apply a learning model, trained by using the samples drawn from a certain domain (*source domain*), to the samples drawn from another domain (*target domain*) with a different distribution (see Bickel *et al.*, 2007; Wu and Dietterich, 2004; Blitzer *et al.*, 2006; Ben-David *et al.*, 2010; Bian *et al.*, 2012).

Without loss of generality, this paper is mainly concerned with two types of representative domain adaptation. In the first type, the learner receives training data from several source domains, known as *domain adaptation with multiple sources* (see Crammer *et al.*, 2006, 2008; Mansour *et al.*, 2008, 2009a). In the second type, the learner minimizes a convex combination of the source and the target empirical risks, termed as *domain adaptation combining source and target data* (see Ben-David *et al.*, 2010; Blitzer *et al.*, 2008).

## 1.1 Overview of Main Results

In this paper, we present a new framework to obtain the generalization bounds of the learning process for the aforementioned two kinds of representative domain adaptation, respectively. Based on the resultant bounds, we then analyze the asymptotical properties of the learning processes for the two types of domain adaptation. There are four major aspects in the framework:

- the quantity measuring the difference between two domains;
- the complexity measure of function class;
- the deviation inequalities of the learning process for domain adaptation;
- the symmetrization inequality of the learning process for domain adaptation.

Generally, in order to obtain the generalization bounds of a learning process, one needs to develop the related deviation (or concentration) inequalities of the learning process. For either kind of domain adaptation, we use a martingale method to develop the related Hoeffding-type deviation inequality. Moreover, in the situation of domain adaptation, since the source domain differs from the target domain, the desired symmetrization inequality for domain adaptation should incorporate some quantity to reflect the difference. From the point of this view, we then obtain the related symmetrization inequality incorporating the integral probability metric that measures the difference between the distributions of the source and the target domains. Next, we present the generalization bounds based on the uniform entropy number for both kinds of domain adaptation. Also, we generalize the classical McDiarmid's inequality to a more general setting, where independent random variables take values from different domains. By using the derived

inequality, we obtain the generalization bounds based on the Rademacher complexity. Following the resultant bounds, we study the asymptotic convergence and the rate of convergence of the learning process in addition to a discussion on factors that affect the asymptotic behaviors. The numerical experiments support our theoretical findings as well. Meanwhile, we give a comparison with the related results under the assumption of same distribution.

## 1.2 Organization of the Paper

The rest of this paper is organized as follows. Section 2 introduces the problems studied in this paper. Section 3 introduces the integral probability metric to measure the difference between two domains. In Section 4, we introduce two kinds of complexity measures of function classes including the uniform entropy number and the Rademacher complexity. In Section 5 (resp. Section 6), we present the generalization bounds of the learning process for domain adaptation with multiple sources (resp. combining source and target data), and then analyze the asymptotic behavior of the learning process in addition to the related numerical experiment supporting our findings. In Section 7, we list the existing works on the theoretical analysis of domain adaptation as a comparison and the last section concludes the paper. In the appendices, we prove main results of this paper. For clarity of presentation, we also postpone the discussion of the deviation inequalities and the symmetrization inequalities in the appendices.

## 2 Problem Setup

In this section, we formalize the main issues of this paper by introducing some necessary notations

### 2.1 Domain Adaptation with Multiple Sources

We denote  $\mathcal{Z}^{(S_k)} := \mathcal{X}^{(S_k)} \times \mathcal{Y}^{(S_k)} \subset \mathbb{R}^I \times \mathbb{R}^J$  ( $1 \leq k \leq K$ ) and  $\mathcal{Z}^{(T)} := \mathcal{X}^{(T)} \times \mathcal{Y}^{(T)} \subset \mathbb{R}^I \times \mathbb{R}^J$  as the  $k$ -th source domain and the target domain, respectively. Set  $L = I + J$ . Let  $\mathcal{D}^{(S_k)}$  and  $\mathcal{D}^{(T)}$  stand for the distributions of the input spaces  $\mathcal{X}^{(S_k)}$  ( $1 \leq k \leq K$ ) and  $\mathcal{X}^{(T)}$ , respectively. Denote  $g_*^{(S_k)} : \mathcal{X}^{(S_k)} \rightarrow \mathcal{Y}^{(S_k)}$  and  $g_*^{(T)} : \mathcal{X}^{(T)} \rightarrow \mathcal{Y}^{(T)}$  as the labeling functions of  $\mathcal{Z}^{(S_k)}$  ( $1 \leq k \leq K$ ) and  $\mathcal{Z}^{(T)}$ , respectively. In the situation of domain adaptation with multiple sources, the input-space distributions  $\mathcal{D}^{(S_k)}$  ( $1 \leq k \leq K$ ) and  $\mathcal{D}^{(T)}$  differ from each other, or  $g_*^{(S_k)}$  ( $1 \leq k \leq K$ ) and  $g_*^{(T)}$  differ from each other, or both of the cases occur. There are sufficient amounts of i.i.d. samples  $\mathbf{Z}_1^{N_k} = \{\mathbf{z}_n^{(k)}\}_{n=1}^{N_k}$  drawn from each source domain  $\mathcal{Z}^{(S_k)}$  ( $1 \leq k \leq K$ ) but little or no labeled samples drawn from the target domain  $\mathcal{Z}^{(T)}$ .

Given  $\mathbf{w} = (w_1, \dots, w_K) \in [0, 1]^K$  with  $\sum_{k=1}^K w_k = 1$ , let  $g_{\mathbf{w}} \in \mathcal{G}$  be the function that minimizes the empirical risk

$$\mathbb{E}_{\mathbf{w}}^{(S)}(\ell \circ g) = \sum_{k=1}^K w_k \mathbb{E}_{N_k}^{(S_k)}(\ell \circ g) = \sum_{k=1}^K \frac{w_k}{N_k} \sum_{n=1}^{N_k} \ell(g(\mathbf{x}_n^{(k)}), \mathbf{y}_n^{(k)}) \quad (1)$$

over  $\mathcal{G}$  with respect to sample sets  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K$ , and it is expected that  $g_{\mathbf{w}}$  will perform well on the target expected risk:

$$\mathbb{E}^{(T)}(\ell \circ g) := \int \ell(g(\mathbf{x}^{(T)}), \mathbf{y}^{(T)}) d\mathbb{P}(\mathbf{z}^{(T)}), \quad g \in \mathcal{G}, \quad (2)$$

i.e.,  $g_{\mathbf{w}}$  approximates the labeling  $g_*^{(T)}$  as precisely as possible.

In the learning process of domain adaptation with multiple sources, we are mainly interested in the following two types of quantities:

- $E^{(T)}(\ell \circ g_{\mathbf{w}}) - E_{\mathbf{w}}^{(S)}(\ell \circ g_{\mathbf{w}})$ , which corresponds to the estimation of the expected risk in the target domain  $\mathcal{Z}^{(T)}$  from the empirical quantity that is the weighted combination of the empirical risks in the multiple sources  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$ ;
- $E^{(T)}(\ell \circ g_{\mathbf{w}}) - E^{(T)}(\ell \circ \tilde{g}_*)$ , which corresponds to the performance of the algorithm for domain adaptation with multiple sources,

where  $\tilde{g}_* \in \mathcal{G}$  is the function that minimizes the expected risk  $E^{(T)}(\ell \circ g)$  over  $\mathcal{G}$ .

Recalling (1) and (2), since

$$E_{\mathbf{w}}^{(S)}(\ell \circ \tilde{g}_*) - E_{\mathbf{w}}^{(S)}(\ell \circ g_{\mathbf{w}}) \geq 0,$$

we have

$$\begin{aligned} E^{(T)}(\ell \circ g_{\mathbf{w}}) &= E^{(T)}(\ell \circ g_{\mathbf{w}}) - E^{(T)}(\ell \circ \tilde{g}_*) + E^{(T)}(\ell \circ \tilde{g}_*) \\ &\leq E_{\mathbf{w}}^{(S)}(\ell \circ \tilde{g}_*) - E_{\mathbf{w}}^{(S)}(\ell \circ g_{\mathbf{w}}) + E^{(T)}(\ell \circ g_{\mathbf{w}}) - E^{(T)}(\ell \circ \tilde{g}_*) + E^{(T)}(\ell \circ \tilde{g}_*) \\ &\leq 2 \sup_{g \in \mathcal{G}} |E^{(T)}(\ell \circ g) - E_{\mathbf{w}}^{(S)}(\ell \circ g)| + E^{(T)}(\ell \circ \tilde{g}_*), \end{aligned} \quad (3)$$

and thus

$$0 \leq E^{(T)}(\ell \circ g_{\mathbf{w}}) - E^{(T)}(\ell \circ \tilde{g}_*) \leq 2 \sup_{g \in \mathcal{G}} |E^{(T)}(\ell \circ g) - E_{\mathbf{w}}^{(S)}(\ell \circ g)|.$$

This shows that the asymptotic behaviors of the aforementioned two quantities when the sample numbers  $N_1, \dots, N_K$  go to *infinity* can both be described by the supremum

$$\sup_{g \in \mathcal{G}} |E^{(T)}(\ell \circ g) - E_{\mathbf{w}}^{(S)}(\ell \circ g)|, \quad (4)$$

which is the so-called generalization bound of the learning process for domain adaptation with multiple sources.

For convenience, we define the loss function as class

$$\mathcal{F} := \{\mathbf{z} \mapsto \ell(g(\mathbf{x}), \mathbf{y}) : g \in \mathcal{G}\}, \quad (5)$$

and call  $\mathcal{F}$  as the function class in the rest of this paper. By (1) and (2), given sample sets  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K$  drawn from  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$  respectively, we briefly denote for any  $f \in \mathcal{F}$ ,

$$E^{(T)}f := \int f(\mathbf{z}^{(T)}) dP(\mathbf{z}^{(T)}), \quad (6)$$

and

$$E_{\mathbf{w}}^{(S)}f := \sum_{k=1}^K \frac{w_k}{N_k} \sum_{n=1}^{N_k} f(\mathbf{z}_n^{(k)}). \quad (7)$$

Thus, we rewrite the generalization bound (4) for domain adaptation with multiple sources as

$$\sup_{f \in \mathcal{F}} |E^{(T)}f - E_{\mathbf{w}}^{(S)}f|. \quad (8)$$

## 2.2 Domain Adaptation Combining Source and Target Data

Denote  $\mathcal{Z}^{(S)} := \mathcal{X}^{(S)} \times \mathcal{Y}^{(S)} \subset \mathbb{R}^I \times \mathbb{R}^J$  and  $\mathcal{Z}^{(T)} := \mathcal{X}^{(T)} \times \mathcal{Y}^{(T)} \subset \mathbb{R}^I \times \mathbb{R}^J$  as the source domain and the target domain, respectively. Let  $\mathcal{D}^{(S)}$  and  $\mathcal{D}^{(T)}$  stand for the distributions of the input spaces  $\mathcal{X}^{(S)}$  and  $\mathcal{X}^{(T)}$ , respectively. Denote  $g_*^{(S)} : \mathcal{X}^{(S)} \rightarrow \mathcal{Y}^{(S)}$  and  $g_*^{(T)} : \mathcal{X}^{(T)} \rightarrow \mathcal{Y}^{(T)}$  as the labeling functions of  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ , respectively. In the situation of domain adaptation combining source and target data (see Blitzer *et al.*, 2008; Ben-David *et al.*, 2010), the input-space distributions  $\mathcal{D}^{(S)}$  and  $\mathcal{D}^{(T)}$  differ from each other, or the labeling functions  $g_*^{(S)}$  and  $g_*^{(T)}$  differ from each other, or both cases occur. There are some (but not enough) samples  $\mathbf{Z}_1^{N_T} := \{\mathbf{z}_n^{(T)}\}_{n=1}^{N_T}$  drawn from the target domain  $\mathcal{Z}^{(T)}$  in addition to a large amount of samples  $\mathbf{Z}_1^{N_S} := \{\mathbf{z}_n^{(S)}\}_{n=1}^{N_S}$  drawn from the source domain  $\mathcal{Z}^{(S)}$  with  $N^{(T)} \ll N^{(S)}$ . Given a  $\tau \in [0, 1)$ , we denote  $g_\tau \in \mathcal{G}$  as the function that minimizes the convex combination of the source and the target empirical risks over  $\mathcal{G}$ :

$$E_\tau(\ell \circ g) := \tau E_{N_T}^{(T)}(\ell \circ g) + (1 - \tau) E_{N_S}^{(S)}(\ell \circ g), \quad (9)$$

and it is expected that  $g_\tau$  will perform well for any pair  $\mathbf{z}^{(T)} = (\mathbf{x}^{(T)}, \mathbf{y}^{(T)}) \in \mathcal{Z}^{(T)}$ , *i.e.*,  $g_\tau$  approximates the labeling function  $g_*^{(T)}$  as precisely as possible.

As mentioned by Blitzer *et al.* (2008); Ben-David *et al.* (2010), setting  $\tau$  involves a tradeoff between the source data that are sufficient but not accurate and the target data that are accurate but not sufficient. Especially, setting  $\tau = 0$  provides a learning process of the basic domain adaptation with one single source (see Ben-David *et al.*, 2006).

Similar to the situation of domain adaptation with multiple sources, two types of quantities:  $E^{(T)}(\ell \circ g_\tau) - E_\tau(\ell \circ g_\tau)$  and  $E^{(T)}(\ell \circ g_\tau) - E^{(T)}(\ell \circ \tilde{g}_*)$  also play an essential role in analyzing the asymptotic behavior of the learning process for domain adaptation combining source and target data. By the similar way of (3), we need to consider the supremum

$$\sup_{g \in \mathcal{G}} |E^{(T)}(\ell \circ g) - E_\tau(\ell \circ g)|, \quad (10)$$

which is the so-called generalization bound of the learning process for domain adaptation combining source and target data. Following the notation of (5) and taking  $f = \ell \circ g$ , we can equivalently rewrite the generalization bound (10) as

$$\sup_{f \in \mathcal{F}} |E^{(T)}f - E_\tau f|. \quad (11)$$

## 3 Integral Probability Metric

As shown in some existing works (see Mansour *et al.*, 2008, 2009a; Ben-David *et al.*, 2010, 2006), one of major challenges in the theoretical analysis of domain adaptation is to find a quantity to measure the difference between the source domain  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$ . Then, one can use the quantity to achieve generalization bounds for domain adaptation. In this section, we use the integral probability metric to measure the difference between the distributions of  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ , and then discuss the relationship between the integral probability metric and other quantities proposed in existing works, *e.g.*, the  $\mathcal{H}$ -divergence and the discrepancy distance (see Ben-David *et al.*, 2010; Mansour *et al.*, 2009b). Moreover, we will show that there is a special situation of domain adaptation, where the integral probability metric performs better than other quantities (see Remark 3.1)

### 3.1 Integral Probability Metric

In Ben-David *et al.* (2010, 2006), the  $\mathcal{H}$ -divergence was introduced to derive the generalization bounds based on the VC dimension under the condition of “ $\lambda$ -close”. Mansour *et al.* (2009b) obtained the generalization bounds based on the Rademacher complexity by using the *discrepancy distance*. Both quantities are aimed to measure the difference between two input-space distributions  $\mathcal{D}^{(S)}$  and  $\mathcal{D}^{(T)}$ . Moreover, Mansour *et al.* (2009a) used the Rényi divergence to measure the distance between two distributions. In this paper, we use the following quantity to measure the difference between the distributions of the source and the target domains:

**Definition 3.1** *Given two domains  $\mathcal{Z}^{(S)}, \mathcal{Z}^{(T)} \subset \mathbb{R}^L$ , let  $\mathbf{z}^{(S)}$  and  $\mathbf{z}^{(T)}$  be the random variables taking values from  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ , respectively. Let  $\mathcal{F} \subset \mathbb{R}^{\mathcal{Z}}$  be a function class. We define*

$$D_{\mathcal{F}}(S, T) := \sup_{f \in \mathcal{F}} |\mathbb{E}^{(S)} f - \mathbb{E}^{(T)} f|, \quad (12)$$

where the expectations  $\mathbb{E}^{(S)}$  and  $\mathbb{E}^{(T)}$  are taken on the distributions  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ , respectively.

The quantity  $D_{\mathcal{F}}(S, T)$  is termed as the integral probability metric that has played an important role in probability theory for measuring the difference between the two probability distributions (see Zolotarev, 1984; Rachev, 1991; Müller, 1997; Reid and Williamson, 2011). Recently, Sriperumbudur *et al.* (2009, 2012) gave the further investigation and proposed an empirical method to compute the integral probability metric. As mentioned by Müller (1997)[page 432], the quantity  $D_{\mathcal{F}}(S, T)$  is a semimetric and it is a metric if and only if the function class  $\mathcal{F}$  separates the set of all signed measures with  $\mu(\mathcal{Z}) = 0$ . Namely, according to Definition 3.1, given a non-trivial function class  $\mathcal{F}$ , the integral probability metric  $D_{\mathcal{F}}(S, T)$  is equal to zero if the domains  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$  have the same distribution.

By (5), the quantity  $D_{\mathcal{F}}(S, T)$  can be equivalently rewritten as

$$\begin{aligned} D_{\mathcal{F}}(S, T) &= \sup_{g \in \mathcal{G}} \left| \mathbb{E}^{(S)} \ell(g(\mathbf{x}^{(S)}), \mathbf{y}^{(S)}) - \mathbb{E}^{(T)} \ell(g(\mathbf{x}^{(T)}), \mathbf{y}^{(T)}) \right| \\ &= \sup_{g \in \mathcal{G}} \left| \mathbb{E}^{(S)} \ell(g(\mathbf{x}^{(S)}), g_*^{(S)}(\mathbf{x}^{(S)})) - \mathbb{E}^{(T)} \ell(g(\mathbf{x}^{(T)}), g_*^{(T)}(\mathbf{x}^{(T)})) \right|. \end{aligned} \quad (13)$$

Next, based on the equivalent form (13), we discuss the relationships between the quantity  $D_{\mathcal{F}}(S, T)$  and other quantities including the  $\mathcal{H}$ -divergence and the *discrepancy distance*.

### 3.2 Relationship with Other Quantities

Before the formal discussion, we briefly introduce the related quantities proposed in existing works (see Ben-David *et al.*, 2010; Mansour *et al.*, 2009b).

#### 3.2.1 $\mathcal{H}$ -Divergence and Discrepancy Distance

In classification tasks, by setting  $\ell$  as the absolute-value loss function ( $\ell(\mathbf{x}, \mathbf{y}) = |\mathbf{x} - \mathbf{y}|$ ), Ben-David *et al.* (2010) introduced a variant of the  $\mathcal{H}$ -divergence:

$$d_{\mathcal{H} \Delta \mathcal{H}}(\mathcal{D}^{(S)}, \mathcal{D}^{(T)}) = \sup_{g_1, g_2 \in \mathcal{H}} \left| \mathbb{E}^{(S)} \ell(g_1(\mathbf{x}^{(S)}), g_2(\mathbf{x}^{(S)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_2(\mathbf{x}^{(T)})) \right| \quad (14)$$

to achieve VC-dimension-based generalization bounds for domain adaptation under the condition of “ $\lambda$ -close”: there exists a  $\lambda > 0$  such that

$$\lambda \geq \inf_{g \in \mathcal{G}} \left\{ \int \ell(g(\mathbf{x}^{(S)}), \mathbf{y}^{(S)}) dP(\mathbf{z}^{(S)}) + \int \ell(g(\mathbf{x}^{(T)}), \mathbf{y}^{(T)}) dP(\mathbf{z}^{(T)}) \right\}.$$

In both of the classification and regression tasks, given a function class  $\mathcal{G}$  and a loss function  $\ell$ , Mansour *et al.* (2009b) defined the *discrepancy distance* as

$$\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)}) = \sup_{g_1, g_2 \in \mathcal{G}} \left| \mathbb{E}^{(S)} \ell(g_1(\mathbf{x}^{(S)}), g_2(\mathbf{x}^{(S)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_2(\mathbf{x}^{(T)})) \right|, \quad (15)$$

and then used this quantity to obtain the generalization bounds based on the Rademacher complexity.

As mentioned by Mansour *et al.* (2009b), the quantities (14) and (15) match in the setting of classification tasks by setting  $\ell$  as the absolute-value loss function, while the usage of (15) does not require the condition of “ $\lambda$ -close” but the usage of (14) does. Recalling Definition 3.1, since there is no limitation on the function class  $\mathcal{F}$ , the integral probability metric  $D_{\mathcal{F}}(S, T)$  can be used in both classification and regression tasks. Therefore, we only consider the relationship between the integral probability metric  $D_{\mathcal{F}}(S, T)$  and the *discrepancy distance*  $\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)})$ .

### 3.2.2 Relationship between $D_{\mathcal{F}}(S, T)$ and $\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)})$

From Definition 3.1 and (13), we can find that the integral probability metric  $D_{\mathcal{F}}(S, T)$  measures the difference between the distributions of the two domains  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ . However, as addressed in Section 2, if a domain  $\mathcal{Z}^{(S)}$  differs from another domain  $\mathcal{Z}^{(T)}$ , there are three possibilities: the input-space distribution  $\mathcal{D}^{(S)}$  differs from  $\mathcal{D}^{(T)}$ , or  $g_*^{(S)}$  differs from  $g_*^{(T)}$ , or both of them occur. Therefore, it is necessary to consider two kinds of differences: the difference between the input-space distributions  $\mathcal{D}^{(S)}$  and  $\mathcal{D}^{(T)}$  and the difference between the labeling functions  $g_*^{(S)}$  and  $g_*^{(T)}$ . Next, we will show that the integral probability metric  $D_{\mathcal{F}}(S, T)$  can be bounded by using two separate quantities that can measure the difference between  $\mathcal{D}^{(S)}$  and  $\mathcal{D}^{(T)}$  and the difference between  $g_*^{(S)}$  and  $g_*^{(T)}$ , respectively.

As shown in (15), the quantity  $\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)})$  actually measures the difference between the input-space distributions  $\mathcal{D}^{(S)}$  and  $\mathcal{D}^{(T)}$ . Moreover, we introduce another quantity to measure the difference between the labeling functions  $g_*^{(S)}$  and  $g_*^{(T)}$ :

**Definition 3.2** Given a loss function  $\ell$  and a function class  $\mathcal{G}$ , we define

$$Q_{\mathcal{G}}^{(T)}(g_*^{(S)}, g_*^{(T)}) := \sup_{g_1 \in \mathcal{G}} \left| \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(T)}(\mathbf{x}^{(T)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(S)}(\mathbf{x}^{(T)})) \right|. \quad (16)$$

Note that if the loss function  $\ell$  and the function class  $\mathcal{G}$  are both non-trivial (*i.e.*,  $\mathcal{F}$  is non-trivial), the quantity  $Q_{\mathcal{G}}^{(T)}(g_*^{(S)}, g_*^{(T)})$  is a (semi)metric between the labeling functions  $g_*^{(S)}$  and  $g_*^{(T)}$ . In fact, it is not hard to verify that  $Q_{\mathcal{G}}^{(T)}(g_*^{(S)}, g_*^{(T)})$  satisfies the triangle inequality and is equal to zero if  $g_*^{(S)}$  and  $g_*^{(T)}$  match.

By combining (13), (15) and (16), we have

$$\begin{aligned}
\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)}) &= \sup_{g_1, g_2 \in \mathcal{G}} \left| \mathbb{E}^{(S)} \ell(g_1(\mathbf{x}^{(S)}), g_2(\mathbf{x}^{(S)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_2(\mathbf{x}^{(T)})) \right| \\
&\geq \sup_{g_1 \in \mathcal{G}} \left| \mathbb{E}^{(S)} \ell(g_1(\mathbf{x}^{(S)}), g_*^{(S)}(\mathbf{x}^{(S)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(S)}(\mathbf{x}^{(T)})) \right| \\
&= \sup_{g_1 \in \mathcal{G}} \left| \mathbb{E}^{(S)} \ell(g_1(\mathbf{x}^{(S)}), g_*^{(S)}(\mathbf{x}^{(S)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(T)}(\mathbf{x}^{(T)})) \right. \\
&\quad \left. + \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(T)}(\mathbf{x}^{(T)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(S)}(\mathbf{x}^{(T)})) \right| \\
&\geq \sup_{g_1 \in \mathcal{G}} \left| \mathbb{E}^{(S)} \ell(g_1(\mathbf{x}^{(S)}), g_*^{(S)}(\mathbf{x}^{(S)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(T)}(\mathbf{x}^{(T)})) \right| \\
&\quad - \sup_{g_1 \in \mathcal{G}} \left| \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(T)}(\mathbf{x}^{(T)})) - \mathbb{E}^{(T)} \ell(g_1(\mathbf{x}^{(T)}), g_*^{(S)}(\mathbf{x}^{(T)})) \right| \\
&= D_{\mathcal{F}}(S, T) - Q_{\mathcal{G}}^{(T)}(g_*^{(S)}, g_*^{(T)}), \tag{17}
\end{aligned}$$

and thus

$$D_{\mathcal{F}}(S, T) \leq \text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)}) + Q_{\mathcal{G}}^{(T)}(g_*^{(S)}, g_*^{(T)}), \tag{18}$$

which implies that the integral probability metric  $D_{\mathcal{F}}(S, T)$  can be bounded by the summation of the *discrepancy distance*  $\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)})$  and the quantity  $Q_{\mathcal{G}}^{(T)}(g_*^{(S)}, g_*^{(T)})$ , which measure the difference between the input-space distributions  $\mathcal{D}^{(S)}$  and  $\mathcal{D}^{(T)}$  and the difference between the labeling functions  $g_*^{(S)}$  and  $g_*^{(T)}$ , respectively.

**Remark 3.1** *Note that there is a specific case in the situation of domain adaptation:  $\mathcal{D}^{(S)}$  differs from  $\mathcal{D}^{(T)}$  and meanwhile  $g_*^{(S)}$  differs from  $g_*^{(T)}$ , while the distribution of the domain  $\mathcal{Z}^{(S)}$  matches with that of the domain  $\mathcal{Z}^{(T)}$ . In this case, the integral probability metric  $D_{\mathcal{F}}(S, T)$  equals to zero, but  $\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)})$  or  $Q_{\mathcal{G}}^{(T)}(g_*^{(S)}, g_*^{(T)})$  neither equals to zero. Therefore, the integral probability metric  $D_{\mathcal{F}}(S, T)$  is more suitable for this case than the discrepancy distance  $\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)})$ .*

## 4 Complexity Measures of Function Classes

Generally, the generalization bound of a certain learning process is achieved by incorporating some complexity measure of the function class, *e.g.*, the covering number, the VC dimension and the Rademacher complexity. In this paper, we are mainly concerned with the uniform entropy number and the Rademacher complexity.

### 4.1 Uniform Entropy Number

The uniform entropy number is derived from the concept of the covering number and we refer to Mendelson (2003) for details. The covering number of a function class  $\mathcal{F}$  is defined as follows:

**Definition 4.1** *Let  $\mathcal{F}$  be a function class and  $d$  be a metric on  $\mathcal{F}$ . For any  $\xi > 0$ , the covering number of  $\mathcal{F}$  at radius  $\xi$  with respect to the metric  $d$ , denoted by  $\mathcal{N}(\mathcal{F}, \xi, d)$  is the minimum size of a cover of radius  $\xi$ .*



In some classical results of statistical learning theory, the covering number is applied by letting  $d$  be the distribution-dependent metric. For example, as shown in Theorem 2.3 of Mendelson (2003), one can set  $d$  as the norm  $\ell_1(\mathbf{Z}_1^N)$  and then derive the generalization bound of the i.i.d. learning process by incorporating the expectation of the covering number, *i.e.*,  $\mathbb{E}\mathcal{N}(\mathcal{F}, \xi, \ell_1(\mathbf{Z}_1^N))$ . However, in the situation of domain adaptation, we only know the information of the source domain, while the expectation  $\mathbb{E}\mathcal{N}(\mathcal{F}, \xi, \ell_1(\mathbf{Z}_1^N))$  is dependent on the distributions of the source and the target domains because  $\mathbf{z} = (\mathbf{x}, \mathbf{y})$ . Therefore, the covering number is no longer suitable for our framework to obtain the generalization bounds for domain adaptation. In contrast, the uniform entropy number is distribution-free and thus we choose it as the complexity measure of function classes to derive the generalization bounds for domain adaptation.

Next, we will consider the uniform entropy number of  $\mathcal{F}$  in the situations of two types of domain adaptation: (i) domain adaptation with multiple sources; (ii) domain adaptation combining source and target data, respectively.

#### 4.1.1 Domain Adaptation with Multiple Sources

For clarity of presentation, we give a useful notation for the following discussion. Let  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K := \{\{\mathbf{z}_n^{(k)}\}_{n=1}^{N_k}\}_{k=1}^K$  be the collection of sample sets drawn from multiple sources  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$ , respectively. Denote  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K := \{\{\mathbf{z}'_n^{(k)}\}_{n=1}^{N_k}\}_{k=1}^K$  as the collection of the ghost sample sets drawn from  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$  such that the ghost sample  $\mathbf{z}'_n^{(k)}$  has the same distribution as  $\mathbf{z}_n^{(k)}$  for any  $1 \leq k \leq K$  and any  $1 \leq n \leq N_k$ . Denote  $\mathbf{Z}_1^{2N_k} := \{\mathbf{Z}_1^{N_k}, \mathbf{Z}_1^{N_k}\}$  for any  $1 \leq k \leq K$ . Moreover, given an  $f \in \mathcal{F}$  and a  $\mathbf{w} = (w_1, \dots, w_K) \in [0, 1]^K$  with  $\sum_{k=1}^K w_k = 1$ , we introduce a variant of the  $\ell_1$  norm:

$$\|f\|_{\ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K)} := \sum_{k=1}^K \frac{w_k}{N_k} \sum_{n=1}^{N_k} |f(\mathbf{z}_n^{(k)})|. \quad (19)$$

It is noteworthy that the variant  $\ell_1^{\mathbf{w}}$  of the  $\ell_1$  norm is still a norm on the functional space, which can be directly verified by using the definition of norm, so we omit it here.

In the situation of domain adaptation with multiple sources, by setting the metric  $d$  as  $\ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K)$ , we then define the uniform entropy number of  $\mathcal{F}$  with respect to the metric  $\ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K)$  as

$$\ln \mathcal{N}_1^{\mathbf{w}}(\mathcal{F}, \xi, 2 \sum_{k=1}^K N_k) := \sup_{\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K} \ln \mathcal{N}(\mathcal{F}, \xi, \ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K)). \quad (20)$$

#### 4.1.2 Domain Adaptation Combining Source and Target Data

In the situation of domain adaptation combining source and target data, we have to introduce another variant of the  $\ell_1$  norm on  $\mathcal{F}$ . Let  $\mathbf{Z}_1^{N_S} = \{\mathbf{z}_n^{(S)}\}_{n=1}^{N_S}$  and  $\bar{\mathbf{Z}}_1^{N_T} = \{\mathbf{z}_n^{(T)}\}_{n=1}^{N_T}$  be two sets of samples drawn from the domains  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ , respectively. Given an  $f \in \mathcal{F}$ , we define for any  $\tau \in [0, 1]$ ,

$$\|f\|_{\ell_1^{\tau}(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T})} := \frac{\tau}{N_T} \sum_{n=1}^{N_T} |f(\mathbf{z}_n^{(T)})| + \frac{1-\tau}{N_S} \sum_{n=1}^{N_S} |f(\mathbf{z}_n^{(S)})|. \quad (21)$$

Note that the variant  $\ell_1^{\tau}$  ( $\tau \in [0, 1]$ ) of the norm  $\ell_1$  is still a norm on the functional space, which can be easily verified by using the definition of norm, so we omit it here.

Moreover, let  $\mathbf{Z}_1^{N_S}$  and  $\overline{\mathbf{Z}}_1^{N_T}$  be the ghost sample sets of  $\mathbf{Z}_1^{N_S}$  and  $\overline{\mathbf{Z}}_1^{N_T}$ , respectively. Denote  $\mathbf{Z}_1^{2N_S} := \{\mathbf{Z}_1^{N_S}, \mathbf{Z}_1^{N_S}\}$  and  $\overline{\mathbf{Z}}_1^{2N_T} := \{\overline{\mathbf{Z}}_1^{N_T}, \overline{\mathbf{Z}}_1^{N_T}\}$ , respectively. Then, the uniform entropy number of  $\mathcal{F}$  with respect to the metric  $\ell_1^r(\mathbf{Z})$  is defined as

$$\ln \mathcal{N}_1^r(\mathcal{F}, \xi, 2(N_S + N_T)) := \sup_{\mathbf{Z}} \ln \mathcal{N}(\mathcal{F}, \xi, \ell_1^r(\mathbf{Z})), \quad (22)$$

where  $\mathbf{Z} := \{\mathbf{Z}_1^{2N_S}, \overline{\mathbf{Z}}_1^{2N_T}\}$ .

## 4.2 Rademacher Complexity

The Rademacher complexity is one of the most frequently used complexity measures of function classes and we refer to Van der Vaart and Wellner (1996); Mendelson (2003) for details.

**Definition 4.2** *Let  $\mathcal{F}$  be a function class and  $\{\mathbf{z}_n\}_{n=1}^N$  be a sample set drawn from  $\mathcal{Z}$ . Denote  $\{\sigma_n\}_{n=1}^N$  be a set of random variables independently taking either value from  $\{-1, 1\}$  with equal probability. Rademacher complexity of  $\mathcal{F}$  is defined as*

$$\mathcal{R}(\mathcal{F}) := \mathbb{E} \sup_{f \in \mathcal{F}} \left\{ \frac{1}{N} \left| \sum_{n=1}^N \sigma_n f(\mathbf{z}_n) \right| \right\} \quad (23)$$

with its empirical version

$$\mathcal{R}_N(\mathcal{F}) := \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \left\{ \frac{1}{N} \left| \sum_{n=1}^N \sigma_n f(\mathbf{z}_n) \right| \right\}, \quad (24)$$

where  $\mathbb{E}$  stands for the expectation taken with respect to all random variables  $\{\mathbf{z}_n\}_{n=1}^N$  and  $\{\sigma_n\}_{n=1}^N$ , and  $\mathbb{E}_\sigma$  stands for the expectation only taken with respect to the random variables  $\{\sigma_n\}_{n=1}^N$ .

## 5 Learning Processes of Domain Adaptation with Multiple Sources

In this section, we present two generalization bounds of the learning process for domain adaptation with multiple sources. They are based on the uniform entropy number and the Rademacher complexity, respectively. By using the derived bounds based on the uniform entropy number, we then analyze the asymptotic convergence and the rate of convergence of the learning process. The numerical experiment supports our theoretical analysis as well.

### 5.1 Generalization Bounds

Based on the uniform entropy number defined in (20), a generalization bound for domain adaptation with multiple sources is presented in the following theorem.

**Theorem 5.1** Assume that  $\mathcal{F}$  is a function class consisting of bounded functions with the range  $[a, b]$ . Let  $\mathbf{w} = (w_1, \dots, w_K) \in [0, 1]^K$  with  $\sum_{k=1}^K w_k = 1$ . Then, given an arbitrary  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ , we have for any  $(\prod_{k=1}^K N_k) \geq \frac{8(b-a)^2}{(\xi')^2}$  and any  $\epsilon > 0$ , with probability at least  $1 - \epsilon$ ,

$$\sup_{f \in \mathcal{F}} |E_{\mathbf{w}}^{(S)} f - E^{(T)} f| \leq D_{\mathcal{F}}^{(\mathbf{w})}(S, T) + \left( \frac{\left( \ln \mathcal{N}_1(\mathcal{F}, \xi'/8, 2 \sum_{k=1}^K N_k) - \ln(\epsilon/8) \right)}{\frac{\left( \prod_{k=1}^K N_k \right)}{32(b-a)^2 \left( \sum_{k=1}^K w_k^2 \left( \prod_{i \neq k} N_i \right) \right)}} \right)^{\frac{1}{2}}, \quad (25)$$

where  $\xi' = \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$  and

$$D_{\mathcal{F}}^{(\mathbf{w})}(S, T) := \sum_{k=1}^K w_k D_{\mathcal{F}}(S_k, T). \quad (26)$$

In the above theorem, we show that the generalization bound  $\sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\mathbf{w}}^{(S)} f|$  can be bounded by the right-hand side of (25). Compared to the classical result under the assumption of same distribution (see Mendelson, 2003, Theorem 2.3 and Definition 2.5): with probability at least  $1 - \epsilon$ ,

$$\sup_{f \in \mathcal{F}} |E_N f - E f| \leq O \left( \left( \frac{\ln \mathcal{N}_1(\mathcal{F}, \xi, N) - \ln(\epsilon/8)}{N} \right)^{\frac{1}{2}} \right) \quad (27)$$

with  $E_N f$  being the empirical risk with respect to the sample set  $\mathbf{Z}_1^N$ , there is a discrepancy quantity  $D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$  that is determined by the two factors: the choice of  $\mathbf{w}$  and the integral probability metrics  $D_{\mathcal{F}}(S_k, T)$  ( $1 \leq k \leq K$ ). The two results will coincide if any source domain and the target domain match, *i.e.*,  $D_{\mathcal{F}}(S_k, T) = 0$  holds for any  $1 \leq k \leq K$ .

In order to prove this result, we develop the specific Hoeffding-type deviation inequality and the symmetrization inequality for domain adaptation with multiple sources, respectively. The detailed proof is arranged in Appendix A. Subsequently, we give another generalization bound based on the Rademacher complexity:

**Theorem 5.2** Assume that  $\mathcal{F}$  is a function class consisting of bounded functions with the range  $[a, b]$ . Let  $\mathbf{w} = (w_1, \dots, w_K) \in [0, 1]^K$  with  $\sum_{k=1}^K w_k = 1$ . Then, we have with probability at least  $1 - \epsilon$ ,

$$\sup_{f \in \mathcal{F}} |E_{\mathbf{w}}^{(S)} f - E^{(T)} f| \leq D_{\mathcal{F}}^{(\mathbf{w})}(S, T) + 2 \sum_{k=1}^K w_k \mathcal{R}^{(k)}(\mathcal{F}) + \sqrt{\sum_{k=1}^K \frac{(b-a)^2 w_k^2 \ln(1/\epsilon)}{2N_k}}, \quad (28)$$

where  $D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$  is defined in (26) and  $\mathcal{R}^{(k)}(\mathcal{F})$  ( $1 \leq k \leq K$ ) are the Rademacher complexities on the source domains  $\mathcal{Z}^{(S_k)}$ , respectively.

Similarly, the derived bound (28) coincides with the related classical result under the assumption of same distribution (see Bousquet *et al.*, 2004, Theorem 5), when any source domain of

$\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$  and the target domain  $\mathcal{Z}^{(T)}$  match, *i.e.*,  $D_{\mathcal{F}}^{(\mathbf{w})}(S, T) = D_{\mathcal{F}}(S_k, T) = 0$  holds for any  $1 \leq k \leq K$ . The proof of this theorem is processed by introducing a generalized version of McDiarmid's inequality which allows independent random variables to take values from different domains (see Appendix C).

Subsequently, based on the derived bound (25), we can analyze the asymptotic behavior of the learning process for domain adaptation with multiple sources.

## 5.2 Asymptotic Convergence

In statistical learning theory, it is well-known that the complexity of function class is one of main factors to the asymptotic convergence of the learning process under the assumption of same distribution (Vapnik, 1999; Van der Vaart and Wellner, 1996; Mendelson, 2003).

From Theorem 5.1, we can directly arrive at the following result showing that the asymptotic convergence of the learning process for domain adaptation with multiple sources is affected by the three aspects: the choice of  $\mathbf{w}$ , the discrepancy quantity  $D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$  and the uniform entropy number  $\ln \mathcal{N}_1^{\mathbf{w}}(\mathcal{F}, \xi/8, 2 \sum_{k=1}^K N_k)$ .

**Theorem 5.3** *Assume that  $\mathcal{F}$  is a function class consisting of the bounded functions with the range  $[a, b]$ . Let  $\mathbf{w} = (w_1, \dots, w_K) \in [0, 1]^K$  with  $\sum_{k=1}^K w_k = 1$ . If the following condition holds:*

$$\lim_{N_1, \dots, N_K \rightarrow +\infty} \frac{\ln \mathcal{N}_1^{\mathbf{w}}(\mathcal{F}, \xi/8, 2 \sum_{k=1}^K N_k)}{\frac{(\prod_{k=1}^K N_k)}{32(b-a)^2 (\sum_{k=1}^K w_k^2 (\prod_{i \neq k} N_i))}} < +\infty, \quad (29)$$

*then we have for any  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ ,*

$$\lim_{N_1, \dots, N_K \rightarrow +\infty} \Pr \left\{ \sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\mathbf{w}}^{(S)} f| > \xi \right\} = 0. \quad (30)$$

As shown in Theorem 5.3, if the choice of  $\mathbf{w} \in [0, 1]^K$  and the uniform entropy number  $\ln \mathcal{N}_1^{\mathbf{w}}(\mathcal{F}, \xi'/8, 2 \sum_{k=1}^K N_k)$  satisfy the condition (29) with  $\sum_{k=1}^K w_k = 1$ , the probability of the event that  $\sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\mathbf{w}}^{(S)} f| > \xi$  will converge to *zero* for any  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ , when the sample numbers  $N_1, \dots, N_K$  of multiple sources go to *infinity*, respectively. This is partially in accordance with the classical result of the asymptotic convergence of the learning process under the assumption of same distribution (*cf.* Theorem 2.3 and Definition 2.5 of [22]): the probability of the event that  $\sup_{f \in \mathcal{F}} |Ef - E_N f| > \xi$  will converge to *zero* for any  $\xi > 0$ , if the uniform entropy number  $\ln \mathcal{N}_1(\mathcal{F}, \xi, N)$  satisfies the following:

$$\lim_{N \rightarrow +\infty} \frac{\ln \mathcal{N}_1(\mathcal{F}, \xi, N)}{N} < +\infty. \quad (31)$$

Note that in the learning process of domain adaptation with multiple sources, the uniform convergence of the empirical risk on the source domains to the expected risk on the target domain may not hold, because the limit (30) does not hold for any  $\xi > 0$  but for any  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ . By contrast, the limit (30) holds for all  $\xi > 0$  in the learning process under the assumption of same distribution, if the condition (31) is satisfied. Again, these two results coincide when any source domain and the target domain match, *i.e.*,  $D_{\mathcal{F}}^{(\mathbf{w})}(S, T) = D_{\mathcal{F}}(S_k, T) = 0$  holds for any  $1 \leq k \leq K$ .

Next, we study the rate of convergence of the learning process for domain adaptation with multiple sources.

### 5.3 Rate of Convergence

Recalling (25), we can find that the rate of convergence is affected by the choice of  $\mathbf{w}$ . According to the Cauchy-Schwarz inequality, setting  $w_k = N_k / \sum_{k=1}^K N_k$  ( $1 \leq k \leq K$ ), we have

$$\max \left\{ \frac{(\prod_{k=1}^K N_k)}{32(b-a)^2 (\sum_{k=1}^K w_k^2 (\prod_{i \neq k} N_i))} \right\} = \frac{N_1 + N_2 + \dots + N_K}{32(b-a)^2}, \quad (32)$$

which minimizes the second term of the right-hand side of (25). Thus, by (25), (27) and (32), we find that the fastest rate of convergence of the learning process is up to  $O(1/\sqrt{N})$  which is the same as the classical result (27) of the learning process under the assumption of same distribution if the discrepancy  $D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$  is ignored.

In addition, the bound (28) based on the Rademacher complexity also implies that the rate of convergence of the learning process is affected by the choice of  $\mathbf{w}$ . Again, according to Cauchy-Schwarz inequality, setting  $w_k = \frac{N_k}{\sum_{k=1}^K N_k}$  ( $1 \leq k \leq K$ ) leads to the fastest rate of convergence:

$$\sqrt{\frac{(b-a)^2 \ln(1/\epsilon)}{2 \sum_{k=1}^K N_k}} = O(1/\sqrt{N}),$$

which is in accordance with the aforementioned analysis. The following numerical experiments support our theoretical findings (see Fig. 1).

### 5.4 Numerical Experiment

We have performed the numerical experiments to verify the theoretic analysis of the asymptotic convergence of the learning processes for domain adaptation with multiple sources. Without loss of generality, we only consider the case of  $K = 2$ , *i.e.*, there are two source domains and one target domain. The experiment data are generated in the following way:

For the target domain  $\mathcal{Z}^{(T)} = \mathcal{X}^{(T)} \times \mathcal{Y}^{(T)} \subset \mathbb{R}^{100} \times \mathbb{R}$ , we consider  $\mathcal{X}^{(T)}$  as a Gaussian distribution  $N(0, 1)$  and draw  $\{\mathbf{x}_n^{(T)}\}_{n=1}^{N_T}$  ( $N_T = 4000$ ) from  $\mathcal{X}^{(T)}$  randomly and independently. Let  $\beta \in \mathbb{R}^{100}$  be a random vector of a Gaussian distribution  $N(1, 5)$ , and let the random vector  $R \in \mathbb{R}^{100}$  be a noise term with  $R \sim N(0, 0.5)$ . For any  $1 \leq n \leq N_T$ , we randomly draw  $\beta$  and  $R$  from  $N(1, 5)$  and  $N(0, 0.01)$  respectively, and then generate  $y_n^{(T)} \in \mathcal{Y}$  as follows:

$$y_n^{(T)} = \langle \mathbf{x}_n^{(T)}, \beta \rangle + R. \quad (33)$$

The derived  $\{(\mathbf{x}_n^{(T)}, y_n^{(T)})\}_{n=1}^{N_T}$  ( $N_T = 4000$ ) are the samples of the target domain  $\mathcal{Z}^{(T)}$  and will be used as the test data.

In the similar way, we generate the sample set  $\{(\mathbf{x}_n^{(1)}, y_n^{(1)})\}_{n=1}^{N_1}$  ( $N_1 = 2000$ ) of the source domain  $\mathcal{Z}^{(S_1)} = \mathcal{X}^{(1)} \times \mathcal{Y}^{(1)} \subset \mathbb{R}^{100} \times \mathbb{R}$ : for any  $1 \leq n \leq N_1$ ,

$$y_n^{(1)} = \langle \mathbf{x}_n^{(1)}, \beta \rangle + R, \quad (34)$$

where  $\mathbf{x}_n^{(1)} \sim N(0.5, 1)$ ,  $\beta \sim N(1, 5)$  and  $R \sim N(0, 0.5)$ .

For the source domain  $\mathcal{Z}^{(S_2)} = \mathcal{X}^{(2)} \times \mathcal{Y}^{(2)} \subset \mathbb{R}^{100} \times \mathbb{R}$ , the samples  $\{(\mathbf{x}_n^{(2)}, y_n^{(2)})\}_{n=1}^{N_2}$  ( $N_2 = 2000$ ) are derived in the following way: for any  $1 \leq n \leq N_2$ ,

$$y_n^{(2)} = \langle \mathbf{x}_n^{(2)}, \beta \rangle + R, \quad (35)$$

where  $\mathbf{x}_n^{(2)} \sim N(2, 5)$ ,  $\beta \sim N(1, 5)$  and  $R \sim N(0, 0.5)$ .

In this experiment, we use the method of Least Square Regression<sup>1</sup> to minimize the empirical risk

$$E_w^{(S)}(\ell \circ g) = \frac{w}{N_1} \sum_{n=1}^{N_1} \ell(g(\mathbf{x}_n^{(1)}), y_n^{(1)}) + \frac{(1-w)}{N_2} \sum_{n=1}^{N_2} \ell(g(\mathbf{x}_n^{(2)}), y_n^{(2)}) \quad (36)$$

for different combination coefficients  $w \in \{0.1, 0.3, 0.5, 0.9\}$  and then compute the discrepancy  $|E_w^{(S)}f - E_{N_T}^{(T)}f|$  for each  $N_1 + N_2$ . The initial  $N_1$  and  $N_2$  both equal to 200. Each test is repeated 30 times and the final result is the average of the 30 results. After each test, we increment  $N_1$  and  $N_2$  by 200 until  $N_1 = N_2 = 2000$ . The experiment results are shown in Fig. 1.

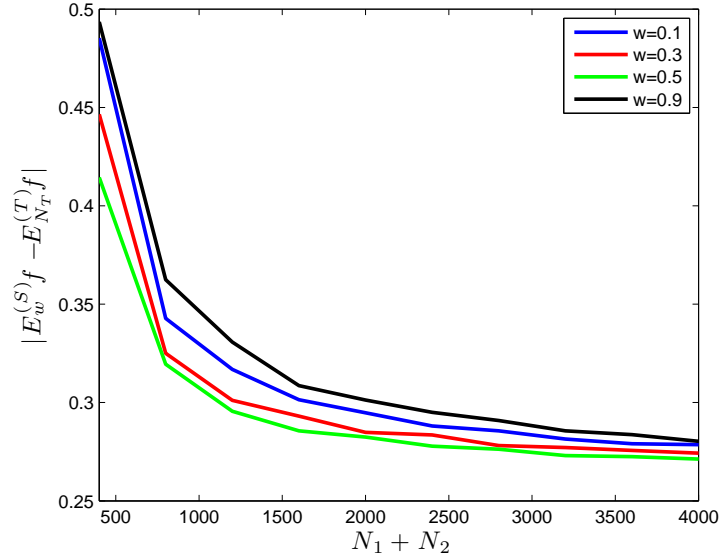


Figure 1: Domain Adaptation with Multiple Sources

From Fig. 1, we can find that for any choice of  $w$ , the curve of  $|E_w^{(S)}f - E_{N_T}^{(T)}f|$  is decreasing when  $N_1 + N_2$  increases, which is in accordance with the results presented in Theorems 5.1 & 5.3. Moreover, when  $w = 0.5$ , the discrepancy  $|E_w^{(S)}f - E_{N_T}^{(T)}f|$  has the fastest rate of convergence, and the rate becomes slower as  $w$  is further away from 0.5. In this experiment, we set  $N_1 = N_2$  that implies that  $N_2/(N_1 + N_2) = 0.5$ . Recalling (25), we have shown that  $w = N_2/(N_1 + N_2)$  will provide the fastest rate of convergence and this proposition is supported by the experiment results shown in Fig. 1.

## 6 Learning Process of Domain Adaptation Combining Source and Target Data

In this section, we present two generalization bounds of the learning process for domain adaptation combining source and target data, which are based on the uniform entropy number and

<sup>1</sup>SLEP Package: <http://www.public.asu.edu/~jye02/Software/SLEP/index.htm>

the Rademacher complexity, respectively. We then analyze the asymptotic convergence and the rate of convergence of the learning process in addition to the numerical experiments supporting our theoretical analysis.

## 6.1 Generalization Bounds

The following theorem provides a generalization bound based on the uniform entropy number with respect to the metric  $\ell_1^\tau$  defined in (22). Similar to the situation of domain adaptation with multiple sources, the proof of this theorem is achieved by using a specific Hoeffding-type deviation inequality and a symmetrization inequality for domain adaptation combining source and target data (see Appendix B).

**Theorem 6.1** *Assume that  $\mathcal{F}$  is a function class consisting of the bounded functions with the range  $[a, b]$ . Let  $\mathbf{Z}_1^{N_S} = \{\mathbf{z}_n^{(S)}\}_{n=1}^{N_S}$  and  $\bar{\mathbf{Z}}_1^{N_T} = \{\mathbf{z}_n^{(T)}\}_{n=1}^{N_T}$  be two sets of i.i.d. samples drawn from domains  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ , respectively. Then, for any  $\tau \in [0, 1)$  and given an arbitrary  $\xi > (1 - \tau)D_{\mathcal{F}}(S, T)$ , we have for any  $N_S N_T \geq \frac{8(b-a)^2}{(\xi')^2}$ , with probability at least  $1 - \epsilon$ ,*

$$\sup_{f \in \mathcal{F}} |E_\tau f - E^{(T)} f| \leq (1 - \tau)D_{\mathcal{F}}(S, T) + \left( \frac{\ln \mathcal{N}_1^\tau(\mathcal{F}, \xi'/8, 2(N_S + N_T)) - \ln(\epsilon/8)}{\frac{N_S N_T}{32(b-a)^2((1-\tau)^2 N_T + \tau^2 N_S)}} \right)^{\frac{1}{2}}, \quad (37)$$

where  $D_{\mathcal{F}}(S, T)$  is defined in (12) and  $\xi' := \xi - (1 - \tau)D_{\mathcal{F}}(S, T)$ .

Compared to the classical result (27) under the assumption of same distribution, the derived bound (37) contains a term of discrepancy quantity  $(1 - \tau)D_{\mathcal{F}}(S, T)$  that is determined by two factors: the combination coefficient  $\tau$  and the quantity  $D_{\mathcal{F}}(S, T)$ . The two results coincide when the source domain  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$  match, i.e.,  $D_{\mathcal{F}}(S, T) = 0$ .

Based on the Rademacher complexity, we then get another generalization bound of the learning process for domain adaptation combining source and target data. Its proof is postponed in Appendix C.

**Theorem 6.2** *Assume that  $\mathcal{F}$  is a function class consisting of the bounded functions with the range  $[a, b]$ . Let  $\mathbf{Z}_1^{N_S} = \{\mathbf{z}_n^{(S)}\}_{n=1}^{N_S}$  and  $\bar{\mathbf{Z}}_1^{N_T} = \{\mathbf{z}_n^{(T)}\}_{n=1}^{N_T}$  be two sets of i.i.d. samples drawn from the domains  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ , respectively. Then, given  $\tau \in [0, 1)$  and for any  $\epsilon > 0$ , we have with probability at least  $1 - \epsilon$ ,*

$$\begin{aligned} \sup_{f \in \mathcal{F}} |E_\tau f - E^{(T)} f| &\leq (1 - \tau)D_{\mathcal{F}}(S, T) + 2(1 - \tau)\mathcal{R}^{(S)}(\mathcal{F}) \\ &\quad + 2\tau\mathcal{R}_{N_T}^{(T)}(\mathcal{F}) + 3\tau\sqrt{\frac{(b-a)\ln(4/\epsilon)}{2N_T}} \\ &\quad + (1 - \tau)\sqrt{\frac{(b-a)^2\ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)}, \end{aligned} \quad (38)$$

where  $D_{\mathcal{F}}(S, T)$  is defined in (12).

Note that in the derived bound (38), we adopt an empirical Rademacher complexity  $\mathcal{R}_{N_T}^{(T)}(\mathcal{F})$  that is based on the data drawn from the target domain  $\mathcal{Z}^{(T)}$ , because the distribution of  $\mathcal{Z}^{(T)}$  is unknown in the situation of domain adaptation. Similar to the aforementioned discussion, the generalization bound (38) coincides with the result under the assumption of same distribution (see Bousquet *et al.*, 2004, Theorem 5), when the source domain of  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$  match, *i.e.*,  $D_{\mathcal{F}}(S, T) = 0$ .

The two results (37) and (38) exhibit a tradeoff between the sample numbers  $N_S$  and  $N_T$ , which is associated with the choice of  $\tau$ . Although the tradeoff has been mentioned in some previous works (see Blitzer *et al.*, 2008; Ben-David *et al.*, 2010), the following will show a rigorous theoretical analysis of the tradeoff.

## 6.2 Asymptotic Convergence

Following Theorem 6.1, we can directly obtain the concerning result pointing out that the asymptotic convergence of the learning process for domain adaptation combining source and target data is affected by three factors: the uniform entropy number  $\ln \mathcal{N}_1^\tau(\mathcal{F}, \xi/8, 2(N_S + N_T))$ , the integral probability metric  $D_{\mathcal{F}}(S, T)$  and the choice of  $\tau \in [0, 1]$ .

**Theorem 6.3** *Assume that  $\mathcal{F}$  is a function class consisting of bounded functions with the range  $[a, b]$ . Given a  $\tau \in [0, 1]$ , if the following condition holds:*

$$\lim_{N_S \rightarrow +\infty} \frac{\ln \mathcal{N}_1^\tau(\mathcal{F}, \xi'/8, 2(N_S + N_T))}{\frac{N_S N_T}{((1-\tau)^2 N_T + \tau^2 N_S)}} < +\infty \quad (39)$$

with  $\xi' := \xi - (1 - \tau)D_{\mathcal{F}}(S, T)$ , then we have for any  $\xi > (1 - \tau)D_{\mathcal{F}}(S, T)$ ,

$$\lim_{N_S \rightarrow +\infty} \Pr \left\{ \sup_{f \in \mathcal{F}} |E^{(T)} f - E_\tau f| > \xi \right\} = 0. \quad (40)$$

As shown in Theorem 6.3, if the choice of  $\tau \in [0, 1]$  and the uniform entropy number  $\ln \mathcal{N}_1^\tau(\mathcal{F}, \xi'/8, 2(N_S + N_T))$  satisfy the condition (39), the probability of the event  $\sup_{f \in \mathcal{F}} |E^{(T)} f - E_\tau f| > \xi$  will converge to zero for any  $\xi > (1 - \tau)D_{\mathcal{F}}(S, T)$ , when  $N_S$  goes to *infinity*. This is partially in accordance with the classical result under the assumption of same distributions derived from the combination of Theorem 2.3 and Definition 2.5 of Mendelson (2003).

Note that in the learning process for domain adaptation combining source and target data, the uniform convergence of the empirical risk  $E_\tau f$  to the expected risk  $E^{(T)} f$  may not hold, because the limit (40) does not hold for any  $\xi > 0$  but for any  $\xi > (1 - \tau)D_{\mathcal{F}}(S, T)$ . By contrast, the limit (40) holds for all  $\xi > 0$  in the learning process under the assumption of same distribution, if the condition (31) is satisfied. The two results coincide when the source domain  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$  match, *i.e.*,  $D_{\mathcal{F}}(S, T) = 0$ .

## 6.3 Rate of Convergence

We consider the choice of  $\tau$  that is an essential factor to the rate of convergence for the learning process and is associated with the tradeoff between the sample numbers  $N_S$  and  $N_T$ . Recalling



(37), if we fix the value of  $\ln \mathcal{N}_1^\tau(\mathcal{F}, \xi'/8, 2(N_S + N_T))$ , setting  $\tau = \frac{N_T}{N_T + N_S}$  minimizes the second term of the right-hand side of (37) and then we arrive at

$$\sup_{f \in \mathcal{F}} |E_\tau f - E^{(T)} f| \leq \frac{N_S D_{\mathcal{F}}(S, T)}{N_S + N_T} + \left( \frac{(\ln \mathcal{N}_1^\tau(\mathcal{F}, \xi'/8, 2(N_S + N_T)) - \ln(\epsilon/8))}{\frac{N_S + N_T}{32(b-a)^2}} \right)^{\frac{1}{2}}, \quad (41)$$

which implies that setting  $\tau = \frac{N_T}{N_T + N_S}$  can result in the fastest rate of convergence, while it can also cause the relatively larger discrepancy between the empirical risk  $E_\tau f$  and the expected risk  $E^{(T)} f$ , because the situation of domain adaptation is set up in the condition that  $N_T \ll N_S$ , which implies that  $\frac{N_S}{N_S + N_T} \approx 1$ . Moreover, this choice of  $\tau$  associated with a trade off between sample numbers  $N_S$  and  $N_T$  is also suitable to the Rademacher-complexity-based bound (38). It is noteworthy that the value  $\tau = \frac{N_T}{N_T + N_S}$  has been mentioned in the section of “Experimental Results” in Blitzer *et al.* (2008). Here, we show a rigorous theoretical analysis of this value and the following numerical experiment also supports this finding (see Fig. 2).

## 6.4 Numerical Experiments

In the situation of domain adaptation combining source and target data, the samples  $\{(\mathbf{x}_n^{(T)}, y_n^{(T)})\}_{n=1}^{N_T}$  ( $N_T = 4000$ ) of the target domain  $\mathcal{Z}^{(T)}$  are generated in the aforementioned way (see (33)). We randomly pick  $N'_T = 100$  samples from them to form the objective function and the rest  $N''_T = 3900$  are used to test.

In the similar way, the samples  $\{(\mathbf{x}_n^{(S)}, y_n^{(S)})\}_{n=1}^{N_S}$  ( $N_S = 4000$ ) of the source domain  $\mathcal{Z}^{(S)}$  are generated as follows: for any  $1 \leq n \leq N_S$ ,

$$y_n^{(S)} = \langle \mathbf{x}_n^{(S)}, \beta \rangle + R, \quad (42)$$

where  $\mathbf{x}_n^{(S)} \sim N(1, 2)$ ,  $\beta \sim N(1, 5)$  and  $R \sim N(0, 0.5)$ .

We also use the method of Least Square Regression to minimize the empirical risk

$$E_\tau(\ell \circ g) = \frac{1 - \tau}{N_S} \sum_{n=1}^{N_S} \ell(g(\mathbf{x}_n^{(1)}), y_n^{(1)}) + \frac{\tau}{N'_T} \sum_{n=1}^{N'_T} \ell(g(\mathbf{x}_n^{(T)}), y_n^{(T)})$$

for different combination coefficients  $\tau \in \{0.1, 0.3, 0.5, 0.9\}$  and then compute the discrepancy  $|E_\tau f - E_{N''_T}^{(T)} f|$  for each  $N_S$ . Since it has to be satisfied that  $N_S \gg N'_T$ , the initial  $N_S$  is set to be 200. Each test is repeated 100 times and the final result is the average of the 100 results. After each test, we increment  $N_S$  by 200 until  $N_S = 4000$ . The experiment results are shown in Fig. 2.

Figure (2) illustrates that for any choice of  $\tau \in \{0.1, 0.3, 0.5, 0.9\}$ , the curve of  $|E_\tau f - E_{N''_T}^{(T)} f|$  is decreasing as  $N_S$  increases. This is in accordance with our results of the asymptotic convergence of the learning process for domain adaptation with multiple sources (see Theorems 6.1 and 6.3). Furthermore, Fig. 2 also shows that when  $\tau \approx N'_T/(N_S + N'_T)$ , the discrepancy  $|E_\tau^{(S)} f - E_{N''_T}^{(T)} f|$  has the fastest rate of convergence, and the rate becomes slower as  $\tau$  is further away from  $N'_T/(N_S + N'_T)$ . Thus, this is in accordance with the theoretical analysis of the asymptotic convergence presented above.

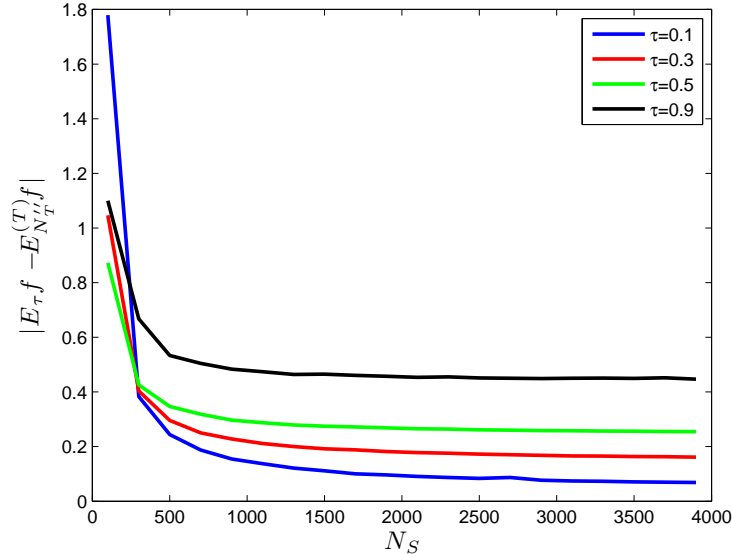


Figure 2: Domain Adaptation Combining Source and Target Data

## 7 Prior Works

There have been some previous works on the theoretical analysis of domain adaptation with multiple sources (see Ben-David *et al.*, 2010; Crammer *et al.*, 2006, 2008; Mansour *et al.*, 2008, 2009a) and domain adaptation combining source and target data (see Blitzer *et al.*, 2008; Ben-David *et al.*, 2010).

In Crammer *et al.* (2006, 2008), the function class and the loss function are assumed to satisfy the conditions of “ $\alpha$ -triangle inequality” and “uniform convergence bound”. Moreover, one has to get some prior information about the disparity between any source domain and the target domain. Under these conditions, some generalization bounds were obtained by using the classical techniques developed under the assumption of same distribution.

Mansour *et al.* (2008) proposed another framework to study the problem of domain adaptation with multiple sources. In this framework, one has to know some prior knowledge including the exact distributions of the source domains and the hypothesis function with a small loss on each source domain. Furthermore, the target domain and the hypothesis function on the target domain were deemed as the mixture of the source domains and the mixture of the hypothesis functions on the source domains, respectively. Then, by introducing the Rényi divergence, Mansour *et al.* (2009a) extended their previous work (Mansour *et al.*, 2008) to a more general setting, where the distribution of the target domain can be arbitrary and one only needs to know an approximation of the exact distribution of each source domain. Ben-David *et al.* (2010) also discussed the situation of domain adaptation with the mixture of source domains.

In Ben-David *et al.* (2010); Blitzer *et al.* (2008), domain adaptation combining source and target data was originally proposed and meanwhile a theoretical framework was presented to analyze its properties for the classification tasks by introducing the  $\mathcal{H}$ -divergence. Under the condition of “ $\lambda$ -close”, the authors applied the classical techniques developed under the assumption of same distribution to achieve the generalization bounds based on the VC dimension.

Mansour *et al.* (2009b) introduced the *discrepancy distance*  $\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)})$  to capture the difference between domains and this quantity can be used in both classification and regression tasks. By applying the classical results of statistical learning theory, the authors obtained the generalization bounds based on the Rademacher complexity.

## 8 Conclusion

In this paper, we propose a new framework to obtain generalization bounds of the learning process for two representative types of domain adaptation: domain adaptation with multiple sources and domain adaptation combining source and target data. This framework is suitable for a variant of learning tasks including classification and regression. Based on the derived bounds, we theoretically analyze the asymptotic convergence and the rate of convergence of the learning process for domain adaptation. There are four important aspects of this framework: the quantity measuring the difference between two domains; the complexity measure of function class, the deviation inequality and the symmetrization inequality for domain adaptation.

- We use the integral probability metric  $D_{\mathcal{F}}(S, T)$  to measure the difference between two domains  $\mathcal{Z}^{(S)}$  and  $\mathcal{Z}^{(T)}$ . We show that the integral probability metric is well-defined and is a (semi)metric on the space of the probability distributions. It can be bounded by the summation of the *discrepancy distance*  $\text{disc}_\ell(\mathcal{D}^{(S)}, \mathcal{D}^{(T)})$  and the quantity  $Q_{\mathcal{G}}^{(T)}(g_*^{(S)}, g_*^{(T)})$ , which measure the difference between the input-space distributions  $\mathcal{D}^{(S)}$  and  $\mathcal{D}^{(T)}$  and the difference between labeling functions  $g_*^{(S)}$  and  $g_*^{(T)}$ , respectively. Note that there is a special case that is more suitable to the integral probability metric  $D_{\mathcal{F}}(S, T)$  than other quantities (see Remark 3.1).
- The uniform entropy number and the Rademacher complexity are adopted to achieved the generalization bounds (25); (37) and (28); (38), respectively. It is noteworthy that the generalization bounds (25) and (37) can lead to the results based on the fat-shattering dimension, respectively (see Mendelson, 2003, Theorem 2.18). According to Theorem 2.6.4 of Van der Vaart and Wellner (1996), the bounds based on the VC dimension can also be obtained from the results (25) and (37), respectively.
- Instead of directly applying the classical techniques, we present the specific deviation inequalities for the learning process of domain adaptation. In order to obtain the generalization bounds based on the uniform entropy numbers, we develop the specific Hoeffding-type deviation inequalities for the two types of domain adaptation, respectively (see Appendices A & B). Furthermore, we also generalize the classical McDiarmid's inequality to a more general setting where the independent random variables can take value from different domains (see Appendix C).
- We also develop the related symmetrization inequalities of the learning process for domain adaptation. The derived inequalities incorporate the discrepancy term that is determined by the difference between the source and the target domains and reflects the learning-transferring from the source to the target domains.

Based on the derived generalization bounds, we provide a rigorous theoretical analysis of the asymptotic convergence and the rate of convergence of the learning process for either kind of

domain adaptation. We also consider the choices of  $\mathbf{w}$  and  $\tau$  that affect the rate of convergence of the learning processes for the two types of domain adaptation, respectively. Moreover, we give a comparison with the previous works Ben-David *et al.* (2010); Crammer *et al.* (2006, 2008); Mansour *et al.* (2008, 2009a); Blitzer *et al.* (2008) as well as the related results of the learning process under the assumption of same distribution (see Bousquet *et al.*, 2004; Mendelson, 2003). The numerical experiments support our theoretical findings as well.

In our future work, we will attempt to find a new distance between distributions to develop the generalization bounds based on other complexity measures, and analyze other theoretical properties of domain adaptation.

## References

- V. N. Vapnik. *The Nature of Statistical Learning Theory (Information Science and Statistics)*. Springer, New York, 1999.
- O. Bousquet, S. Boucheron and G. Lugosi. Introduction to Statistical Learning Theory. *Lecture Notes in Artificial Intelligence*, 3176:169-207, 2004.
- A. van der Vaart and J. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics (Springer Series in Statistics)*. Springer, 1996.
- A. Blumer, A. Ehrenfeucht, D. Haussler, and M.K. Warmuth. Learnability and the Vapnik-Chervonenkis Dimension. *Journal of the ACM*, 36(4):929-965, 1989.
- P.L. Bartlett, O. Bousquet and S. Mendelson. Local Rademacher Complexities. *Annals of Statistics*, 33:1497-1537, 2005.
- Z. Hussain and J. Shawe-Taylor. Improved Loss Bounds for Multiple Kernel Learning. *Journal of Machine Learning Research*, **15(W&CP)**:370-377, 2011.
- J. Jiang, and C. Zhai. Instance Weighting for Domain Adaptation in NLP. *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics (ACL)*, 264-271, 2007.
- J. Blitzer, M. Dredze and F. Pereira. Biographies, bollywood, boomboxes and blenders: Domain adaptation for sentiment classification. *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics (ACL)*, 440-447, 2007.
- S. Bickel, M. Brückner and T. Scheffer. Discriminative learning for differing training and test distributions. *Proceedings of the 24th international conference on Machine learning (ICML)*, 81-88, 2007.
- P. Wu and T.G. Dietterich. Improving SVM accuracy by training on auxiliary data sources. *Proceedings of the twenty-first international conference on Machine learning*, 2004.
- J. Blitzer, R. McDonald and F. Pereira. Domain adaptation with structural correspondence learning. *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2006.
- S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira and J. Wortman. A Theory of Learning from Different Domains. *Machine Learning*, 79:151-175, 2010.

- W. Bian, D. Tao and Y. Rui. Cross-Domain Human Action Recognition. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 42(2):298-307, 2012.
- K. Crammer, M. Kearns and J. Wortman. Learning from Multiple Sources. *Advances in Neural Information Processing Systems (NIPS)*, 2006.
- K. Crammer, M. Kearns and J. Wortman. Learning from Multiple Sources. *Journal of Machine Learning Research*, 9:1757-1774, 2008.
- Y. Mansour, M. Mohri and A. Rostamizadeh. Domain adaptation with multiple sources. *Advances in Neural Information Processing Systems (NIPS)*, 2008.
- Y. Mansour, M. Mohri and A. Rostamizadeh. Multiple Source Adaptation and The Rényi Divergence. *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence (UAI)*, 2000.
- J. Blitzer, K. Crammer, A. Kulesza, F. Pereira and J. Wortman. Learning Bounds for Domain Adaptation. *Advances in Neural Information Processing Systems (NIPS)*, 2008.
- Y. Mansour, M. Mohri and A. Rostamizadeh. Domain Adaptation: Learning Bounds and Algorithms. *Conference on Learning Theory (COLT)*, 2009.
- S. Ben-David, J. Blitzer, K. Crammer and F. Pereira. Analysis of Representations for Domain Adaptation. *Advances in Neural Information Processing Systems (NIPS)*, 2006.
- V.M. Zolotarev. Probability Metrics. *Theory of Probability and its Application*, 28(1):278-302, 1984.
- S.T. Rachev. *Probability Metrics and the Stability of Stochastic Models*. John Wiley and Sons, 1991.
- A. Müller. Integral Probability Metrics and Their Generating Classes of Functions. *Advances in Applied Probability*, 29(2):429-443, 1997.
- M.D. Reid and R.C. Williamson. Information, Divergence and Risk for Binary Experiments. *Journal of Machine Learning Research*, 12:731-817, 2011.
- B.K. Sriperumbudur, A. Gretton, K. Fukumizu, G.R.G. Lanckriet and B. Schölkopf. A Note on Integral Probability Metrics and  $\phi$ -Divergences. *CoRR*, abs/0901.2698, 2009.
- B.K. Sriperumbudur, K. Fukumizu, A. Gretton, B. Schölkopf and G.R.G. Lanckriet. On the empirical estimation of integral probability metrics. *Electron. J. Statist.*, 6:1550-1599, 2012.
- S. Mendelson. A Few Notes on Statistical Learning Theory. *Advanced Lectures on Machine Learning*, 2600:1-40, 2003.
- W. Hoeffding. Probability Inequalities for Sums of Bounded Random Variables. *Journal of the American Statistical Association*, 58(301):13-30, 1963.

## A Proof of Theorem 5.1

In this appendix, we provide the proof of Theorem 5.1. In order to achieve the proof, we need to develop the specific Hoeffding-type deviation inequality and the symmetrization inequality for domain adaptation with multiple sources.

### A.1 Hoeffding-Type Deviation Inequality for Multiple Sources

Deviation (or concentration) inequalities play an essential role in obtaining the generalization bounds for a certain learning process. Generally, specific deviation inequalities need to be developed for different learning processes. There are many popular deviation and concentration inequalities, *e.g.*, Hoeffding's inequality, McDiarmid's inequality, Bennett's inequality, Bernstein's inequality and Talagrand's inequality. These results are all built under the assumption of same distribution, and thus they are not applicable (or at least cannot be directly applied) to the setting of multiple sources. Next, based on Hoeffding's inequality (Hoeffding, 1963), we present a deviation inequality for multiple sources.

**Theorem A.1** *Assume that  $\mathcal{F}$  is a function class consisting of bounded functions with the range  $[a, b]$ . Let  $\mathbf{Z}_1^{N_k} = \{\mathbf{z}_n^{(k)}\}_{n=1}^{N_k}$  be the set of i.i.d. samples drawn from the source domain  $\mathcal{Z}^{(S_k)} \subset \mathbb{R}^L$  ( $1 \leq k \leq K$ ). Given  $\mathbf{w} \in [0, 1]^K$  with  $\sum_{k=1}^K w_k = 1$  and for any  $f \in \mathcal{F}$ , we define a function  $F_{\mathbf{w}} : \mathbb{R}^{L \sum_{k=1}^K N_k} \rightarrow \mathbb{R}$  as*

$$F_{\mathbf{w}} \left( \{\mathbf{X}_1^{N_k}\}_{k=1}^K \right) := \sum_{k=1}^K w_k \left( \prod_{i \neq k} N_i \right) \sum_{n=1}^{N_k} f(\mathbf{x}_n^{(k)}), \quad (43)$$

where for any  $1 \leq k \leq K$  and given  $N_k \in \mathbb{N}$ , the set  $\mathbf{X}_1^{N_k}$  is denoted as

$$\mathbf{X}_1^{N_k} := \{\mathbf{x}_1^{(k)}, \mathbf{x}_2^{(k)}, \dots, \mathbf{x}_{N_k}^{(k)}\} \in (\mathbb{R}^L)^{N_k}.$$

Then, we have for any  $\xi > 0$ ,

$$\begin{aligned} & \Pr \left\{ \left| \mathbb{E}^{(S)} F_{\mathbf{w}} - F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) \right| > \xi \right\} \\ & \leq 2 \exp \left\{ - \frac{2\xi^2}{(b-a)^2 \left( \prod_{k=1}^K N_k \right) \left( \sum_{k=1}^K w_k^2 \left( \prod_{i \neq k} N_i \right) \right)} \right\}, \end{aligned} \quad (44)$$

where  $\mathbb{E}^{(S)}$  stands for the expectation taken on all source domains  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$ .

This result is an extension of the classical Hoeffding-type deviation inequality under the assumption of same distribution (see Bousquet *et al.*, 2004, Theorem 1). Compared to the classical result, the resultant deviation inequality (44) is suitable to the setting of multiple sources. These two inequalities coincide when there is only one source, *i.e.*,  $K = 1$ .

The proof of Theorem A.1 is processed by a martingale method. Before the formal proof, we introduce some essential notations.

Let  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K$  be sample sets drawn from multiple sources  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$ , respectively. Define a random variable

$$S_n^{(k)} := \mathbb{E}^{(S)} \left\{ F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) \middle| \mathbf{Z}_1^{N_1}, \mathbf{Z}_1^{N_2}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{Z}_1^n \right\}, \quad 1 \leq k \leq K, \quad 0 \leq n \leq N_k, \quad (45)$$

where

$$\mathbf{Z}_1^n = \{\mathbf{z}_1^{(k)}, \mathbf{z}_2^{(k)}, \dots, \mathbf{z}_n^{(k)}\} \subseteq \mathbf{Z}_1^{N_k}, \text{ and } \mathbf{Z}_1^0 = \emptyset.$$

It is clear that

$$S_0^{(1)} = \mathbb{E}^{(S)} F_{\mathbf{w}} \text{ and } S_{N_K}^{(K)} = F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K),$$

where  $\mathbb{E}^{(S)}$  stands for the expectation taken on all source domains  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$ .

Then, according to (43) and (45), we have for any  $1 \leq k \leq K$  and  $1 \leq n \leq N_k$ ,

$$\begin{aligned} S_n^{(k)} - S_{n-1}^{(k)} &= \mathbb{E}^{(S)} \left\{ F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) \middle| \mathbf{Z}_1^{N_1}, \mathbf{Z}_1^{N_2}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{Z}_1^n \right\} \\ &\quad - \mathbb{E}^{(S)} \left\{ F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) \middle| \mathbf{Z}_1^{N_1}, \mathbf{Z}_1^{N_2}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{Z}_1^{n-1} \right\} \\ &= \mathbb{E}^{(S)} \left\{ \sum_{k=1}^K w_k \left( \prod_{i \neq k} N_i \right) \sum_{n=1}^{N_k} f(\mathbf{z}_n^{(k)}) \middle| \mathbf{Z}_1^{N_1}, \mathbf{Z}_1^{N_2}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{Z}_1^n \right\} \\ &\quad - \mathbb{E}^{(S)} \left\{ \sum_{k=1}^K w_k \left( \prod_{i \neq k} N_i \right) \sum_{n=1}^{N_k} f(\mathbf{z}_n^{(k)}) \middle| \mathbf{Z}_1^{N_1}, \mathbf{Z}_1^{N_2}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{Z}_1^{n-1} \right\} \\ &= \sum_{l=1}^{k-1} w_l \left( \prod_{i \neq l} N_i \right) \sum_{j=1}^{N_l} f(\mathbf{z}_j^{(l)}) + w_k \left( \prod_{i \neq k} N_i \right) \sum_{j=1}^n f(\mathbf{z}_j^{(k)}) \\ &\quad + \mathbb{E}^{(S)} \left\{ \sum_{l=k+1}^K w_l \left( \prod_{i \neq l} N_i \right) \sum_{j=1}^{N_l} f(\mathbf{z}_j^{(l)}) + w_k \left( \prod_{i \neq k} N_i \right) \sum_{j=n+1}^{N_k} f(\mathbf{z}_j^{(k)}) \right\} \\ &\quad - \sum_{l=1}^{k-1} w_l \left( \prod_{i \neq l} N_i \right) \sum_{j=1}^{N_l} f(\mathbf{z}_j^{(l)}) - w_k \left( \prod_{i \neq k} N_i \right) \sum_{j=1}^{n-1} f(\mathbf{z}_j^{(k)}) \\ &\quad - \mathbb{E}^{(S)} \left\{ \sum_{l=k+1}^K w_l \left( \prod_{i \neq l} N_i \right) \sum_{j=1}^{N_l} f(\mathbf{z}_j^{(l)}) + w_k \left( \prod_{i \neq k} N_i \right) \sum_{j=n}^{N_k} f(\mathbf{z}_j^{(k)}) \right\} \\ &= w_k \left( \prod_{i \neq k} N_i \right) (f(\mathbf{z}_n^{(k)}) - \mathbb{E}^{(S_k)} f). \end{aligned} \tag{46}$$

To prove Theorem A.1, we need the following inequality resulted from Hoeffding's lemma.

**Lemma A.1** *Let  $f$  be a function with the range  $[a, b]$ . Then, the following holds for any  $\alpha > 0$ :*

$$\mathbb{E} \left\{ e^{\alpha(f(\mathbf{z}^{(S)}) - \mathbb{E}^{(S)} f)} \right\} \leq e^{\frac{\alpha^2(b-a)^2}{8}}. \tag{47}$$

*Proof.* We consider

$$(f(\mathbf{z}^{(S)}) - \mathbb{E}^{(S)} f)$$

as a random variable. Then, it is clear that

$$\mathbb{E}\{f(\mathbf{z}^{(S)}) - \mathbb{E}^{(S)} f\} = 0.$$

Since the value of  $\mathbb{E}^{(S)} f$  is a constant denoted as  $e$ , we have

$$a - e \leq f(\mathbf{z}^{(S)}) - \mathbb{E}^{(S)} f \leq b - e.$$

According to Hoeffding's lemma, we then have

$$\mathbb{E} \left\{ e^{\alpha(f(\mathbf{z}^{(S)}) - \mathbb{E}^{(S)} f)} \right\} \leq e^{\frac{\alpha^2(b-a)^2}{8}}. \quad (48)$$

This completes the proof.  $\blacksquare$

We are now ready to prove Theorem A.1.

*Proof of Theorem A.1.* According to (43), (46), Lemma A.1, Markov's inequality, Jensen's inequality and the law of iterated expectation, we have for any  $\alpha > 0$ ,

$$\begin{aligned} & \Pr \left\{ F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) - \mathbb{E}^{(S)} F_{\mathbf{w}} > \xi \right\} \\ & \leq e^{-\alpha\xi} \mathbb{E} \left\{ e^{\alpha \left( F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) - \mathbb{E}^{(S)} F_{\mathbf{w}} \right)} \right\} \\ & = e^{-\alpha\xi} \mathbb{E} \left\{ \mathbb{E} \left\{ e^{\alpha \sum_{k=1}^K \sum_{n=1}^{N_k} (S_n^{(k)} - S_{n-1}^{(k)})} \middle| \mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_{K-1}}, \mathbf{Z}_1^{N_K-1} \right\} \right\} \\ & = e^{-\alpha\xi} \mathbb{E} \left\{ e^{\alpha \left( \sum_{k=1}^K \sum_{n=1}^{N_k} (S_n^{(k)} - S_{n-1}^{(k)}) - (S_{N_K}^{(K)} - S_{N_K-1}^{(K)}) \right)} \mathbb{E} \left\{ e^{\alpha (S_{N_K}^{(K)} - S_{N_K-1}^{(K)})} \middle| \mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_{K-1}}, \mathbf{Z}_1^{N_K-1} \right\} \right\} \\ & = e^{-\alpha\xi} \mathbb{E} \left\{ e^{\alpha \left( \sum_{k=1}^K \sum_{n=1}^{N_k} (S_n^{(k)} - S_{n-1}^{(k)}) - (S_{N_K}^{(K)} - S_{N_K-1}^{(K)}) \right)} \mathbb{E} \left\{ e^{\alpha w_K (\prod_{i \neq K} N_i) (f(\mathbf{z}_N^{(K)}) - \mathbb{E}^{(S_K)} f)} \right\} \right\} \\ & \leq e^{-\alpha\xi} \mathbb{E} \left\{ e^{\alpha \left( \sum_{k=1}^K \sum_{n=1}^{N_k} (S_n^{(k)} - S_{n-1}^{(k)}) - (S_{N_K}^{(K)} - S_{N_K-1}^{(K)}) \right)} \right\} e^{\frac{\alpha^2 w_K^2 (\prod_{i \neq K} N_i)^2 (b-a)^2}{8}}, \end{aligned} \quad (49)$$

where  $\mathbf{Z}_1^{N_K-1} := \{\mathbf{z}_1^{(K)}, \dots, \mathbf{z}_{N_K-1}^{(K)}\} \subset \mathbf{Z}_1^{N_K}$ . Therefore, we have

$$\Pr \left\{ F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) - \mathbb{E}^{(S)} F_{\mathbf{w}} > \xi \right\} \leq e^{\Phi(\alpha) - \alpha\xi}, \quad (50)$$

where

$$\Phi(\alpha) = \frac{\alpha^2(b-a)^2 \left( \prod_{k=1}^K N_k \right) \left( \sum_{k=1}^K w_k^2 \left( \prod_{i \neq k} N_i \right) \right)}{8}. \quad (51)$$

Similarly, we can obtain

$$\Pr \left\{ \mathbb{E}^{(S)} F_{\mathbf{w}} - F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) > \xi \right\} \leq e^{\Phi(\alpha) - \alpha\xi}. \quad (52)$$

Note that  $\Phi(\alpha) - \alpha\xi$  is a quadratic function with respect to  $\alpha > 0$  and thus the minimum value “ $\min_{\alpha>0} \{\Phi(\alpha) - \alpha\xi\}$ ” is achieved when

$$\alpha = \frac{4\xi}{(b-a)^2 \left( \prod_{k=1}^K N_k \right) \left( \sum_{k=1}^K w_k^2 \left( \prod_{i \neq k} N_i \right) \right)}.$$

By combining (50), (51) and (52), we arrive at

$$\begin{aligned} & \Pr \left\{ |\mathbb{E}^{(S)} F_{\mathbf{w}} - F_{\mathbf{w}}(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K)| > \xi \right\} \\ & \leq 2 \exp \left\{ - \frac{2\xi^2}{(b-a)^2 \left( \prod_{k=1}^K N_k \right) \left( \sum_{k=1}^K w_k^2 \left( \prod_{i \neq k} N_i \right) \right)} \right\}. \end{aligned}$$



This completes the proof. ■

In the following subsection, we present a symmetrization inequality for domain adaptation with multiple sources.

## A.2 Symmetrization Inequality

Symmetrization inequalities are mainly used to replace the expected risk by an empirical risk computed on another sample set that is independent of the given sample set but has the same distribution. In this manner, the generalization bounds can be achieved by applying some kinds of complexity measures, *e.g.*, the covering number and the VC dimension. However, the classical symmetrization results are built under the assumption of same distribution (see Bousquet *et al.*, 2004). The symmetrization inequality for domain adaptation with multiple sources is presented in the following theorem:

**Theorem A.2** *Assume that  $\mathcal{F}$  is a function class with the range  $[a, b]$ . Let sample sets  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K$  and  $\{\mathbf{Z}'_1^{N_k}\}_{k=1}^K$  be drawn from the source domains  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$ . Then, given an arbitrary  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$  and  $\mathbf{w} = (w_1, \dots, w_K) \in [0, 1]^K$  with  $\sum_{k=1}^K w_k = 1$ , we have for any  $(\prod_{k=1}^K N_k) \geq \frac{8(b-a)^2}{(\xi')^2}$ ,*

$$\Pr \left\{ \sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\mathbf{w}}^{(S)} f| > \xi \right\} \leq 2 \Pr \left\{ \sup_{f \in \mathcal{F}} |E'_{\mathbf{w}}^{(S)} f - E_{\mathbf{w}}^{(S)} f| > \frac{\xi'}{2} \right\}, \quad (53)$$

where  $\xi' = \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ .

This theorem shows that given  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ , the probability of the event:

$$\sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\mathbf{w}}^{(S)} f| > \xi$$

can be bounded by using the probability of the event:

$$\sup_{f \in \mathcal{F}} |E'_{\mathbf{w}}^{(S)} f - E_{\mathbf{w}}^{(S)} f| > \frac{\xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T)}{2} \quad (54)$$

that is only determined by the characteristics of the source domains  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$  when  $\prod_{k=1}^K N_k \geq \frac{8(b-a)^2}{(\xi')^2}$  with  $\xi' = \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ . Compared to the classical symmetrization result under the assumption of same distribution (see Bousquet *et al.*, 2004), there is a discrepancy term  $D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$  in the derived inequality. Especially, the two results coincide when any source domain and the target domain match, *i.e.*,  $D_{\mathcal{F}}(S_k, T) = 0$  holds for any  $1 \leq k \leq K$ . The following is the proof of Theorem A.2.

*Proof of Theorem A.2.* Let  $\hat{f}$  be the function achieving the supremum:

$$\sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\mathbf{w}}^{(S)} f|$$

with respect to the sample set  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K$ . According to (6), (7), (12) and (26), we arrive at

$$|E^{(T)} \hat{f} - E_{\mathbf{w}}^{(S)} \hat{f}| = |E^{(T)} \hat{f} - \bar{E}^{(S)} \hat{f} + \bar{E}^{(S)} \hat{f} - E_{\mathbf{w}}^{(S)} \hat{f}| \leq D_{\mathcal{F}}^{(\mathbf{w})}(S, T) + |\bar{E}^{(S)} \hat{f} - E_{\mathbf{w}}^{(S)} \hat{f}|, \quad (55)$$

and thus,

$$\Pr \left\{ |\mathbb{E}^{(T)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| > \xi \right\} \leq \Pr \left\{ D_{\mathcal{F}}^{(\mathbf{w})}(S, T) + |\bar{\mathbb{E}}^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| > \xi \right\}, \quad (56)$$

where the expectation  $\mathbb{E}^{(S)} \hat{f}$  is defined as

$$\bar{\mathbb{E}}^{(S)} \hat{f} := \sum_{k=1}^K w_k \mathbb{E}^{(S_k)} \hat{f}. \quad (57)$$

Let

$$\xi' := \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T), \quad (58)$$

and denote  $\wedge$  as the conjunction of two events. According to the triangle inequality, we have

$$\left( |\bar{\mathbb{E}}^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| - |\mathbb{E}'^{(S)} \hat{f} - \bar{\mathbb{E}}^{(S)} \hat{f}| \right) \leq |\mathbb{E}'^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}|,$$

and thus for any  $\xi' > 0$ ,

$$\begin{aligned} & \left( \mathbf{1}_{|\bar{\mathbb{E}}^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| > \xi'} \right) \left( \mathbf{1}_{|\mathbb{E}'^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| < \frac{\xi'}{2}} \right) \\ &= \mathbf{1}_{\left\{ |\bar{\mathbb{E}}^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| > \xi' \right\} \wedge \left\{ |\mathbb{E}'^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| < \frac{\xi'}{2} \right\}} \\ &\leq \mathbf{1}_{|\mathbb{E}'^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| > \frac{\xi'}{2}}. \end{aligned}$$

Then, taking the expectation with respect to  $\{\mathbf{Z}'_1^{N_k}\}_{k=1}^K$  gives

$$\begin{aligned} & \left( \mathbf{1}_{|\bar{\mathbb{E}}^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| > \xi'} \right) \Pr' \left\{ |\bar{\mathbb{E}}^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| < \frac{\xi'}{2} \right\} \\ &\leq \Pr' \left\{ |\mathbb{E}'^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| > \frac{\xi'}{2} \right\}. \end{aligned} \quad (59)$$

By Chebyshev's inequality, since  $\{\mathbf{Z}'_1^{N_k}\}_{k=1}^K$  are the sets of i.i.d. samples drawn from the multiple sources  $\{\mathcal{Z}^{(S_k)}\}_{k=1}^K$  respectively, we have for any  $\xi' > 0$ ,

$$\begin{aligned} \Pr' \left\{ |\bar{\mathbb{E}}^{(S)} \hat{f} - \mathbb{E}_{\mathbf{w}}^{(S)} \hat{f}| \geq \frac{\xi'}{2} \right\} &\leq \Pr' \left\{ \sum_{k=1}^K \frac{w_k}{N_k} \sum_{n=1}^{N_k} |\mathbb{E}^{(S_k)} \hat{f} - \hat{f}(\mathbf{z}'_n^{(k)})| \geq \frac{\xi'}{2} \right\} \\ &= \Pr' \left\{ \sum_{k=1}^K w_k \left( \prod_{i \neq k} N_i \right) \sum_{n=1}^{N_k} |\mathbb{E}^{(S_k)} \hat{f} - \hat{f}(\mathbf{z}'_n^{(k)})| \geq \frac{\xi' \prod_{k=1}^K N_k}{2} \right\} \\ &\leq \frac{4\mathbb{E} \left\{ \sum_{k=1}^K w_k \left( \prod_{i \neq k} N_i \right) \sum_{n=1}^{N_k} |\mathbb{E}^{(S_k)} \hat{f} - \hat{f}(\mathbf{z}'_n^{(k)})|^2 \right\}}{\left( \prod_{k=1}^K N_k \right)^2 (\xi')^2} \\ &= \frac{4\mathbb{E} \left\{ \sum_{k=1}^K w_k \left( \prod_{i \neq k} N_i \right) N_k (b-a)^2 \right\}}{\left( \prod_{k=1}^K N_k \right)^2 (\xi')^2} \\ &= \frac{4 \left( \prod_{k=1}^K N_k \right) (b-a)^2}{\left( \prod_{k=1}^K N_k \right)^2 (\xi')^2} = \frac{4(b-a)^2}{(\xi')^2 \left( \prod_{k=1}^K N_k \right)}. \end{aligned} \quad (60)$$

Subsequently, according to (59) and (60), we have for any  $\xi' > 0$ ,

$$\Pr' \left\{ |E_{\mathbf{w}}^{(S)} \hat{f} - E_{\mathbf{w}}^{(S)} \hat{f}| > \frac{\xi'}{2} \right\} \geq \left( \mathbf{1}_{|\bar{E}^{(S)} \hat{f} - E_{\mathbf{w}}^{(S)} \hat{f}| > \xi'} \right) \left( 1 - \frac{4(b-a)^2}{(\xi')^2 (\prod_{k=1}^K N_k)} \right). \quad (61)$$

By combining (56), (58) and (61), taking the expectation with respect to  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K$  and letting

$$\frac{4(b-a)^2}{(\xi')^2 (\prod_{k=1}^K N_k)} \leq \frac{1}{2}$$

can lead to: for any  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ ,

$$\begin{aligned} \Pr \left\{ |E^{(T)} \hat{f} - E_{\mathbf{w}}^{(S)} \hat{f}| > \xi \right\} &\leq \Pr \left\{ |\bar{E}^{(S)} \hat{f} - E_{\mathbf{w}}^{(S)} \hat{f}| > \xi' \right\} \\ &\leq 2\Pr \left\{ |E_{\mathbf{w}}^{(S)} \hat{f} - E_{\mathbf{w}}^{(S)} \hat{f}| > \frac{\xi'}{2} \right\} \end{aligned} \quad (62)$$

with  $\xi' = \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ . This completes the proof.  $\blacksquare$

By using the resultant deviation inequality and the symmetrization inequality, we can achieve the proof of Theorem 5.1.

### A.3 Proof of Theorem 5.1

*Proof of Theorem 5.1.* Consider  $\epsilon$  as an independent Rademacher random variables, *i.e.*, an independent  $\{-1, 1\}$ -valued random variable with equal probability of taking either value. Given sample sets  $\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K$ , denote for any  $f \in \mathcal{F}$  and  $1 \leq k \leq K$ ,

$$\vec{\epsilon}^{(k)} := (\epsilon_1^{(k)}, \dots, \epsilon_{N_k}^{(k)}, -\epsilon_1^{(k)}, \dots, -\epsilon_{N_k}^{(k)}) \in \{-1, 1\}^{2N_k}, \quad (63)$$

and for any  $f \in \mathcal{F}$ ,

$$\vec{f}(\mathbf{Z}_1^{2N_k}) := (f(\mathbf{z}'^{(k)}), \dots, f(\mathbf{z}'_{N_k}^{(k)}), f(\mathbf{z}_1^{(k)}), \dots, f(\mathbf{z}_{N_k}^{(k)})). \quad (64)$$

According to (6), (7) and Theorem A.2, given an arbitrary  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ , we have for any  $\{N_k\}_{k=1}^K \in \mathbb{N}^K$  such that  $\prod_{k=1}^K N_k \geq \frac{8(b-a)^2}{(\xi')^2}$  with  $\xi' = \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ ,

$$\begin{aligned} &\Pr \left\{ \sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\mathbf{w}}^{(S)} f| > \xi \right\} \\ &\leq 2\Pr \left\{ \sup_{f \in \mathcal{F}} |E_{\mathbf{w}}^{(S)} f - E_{\mathbf{w}}^{(S)} f| > \frac{\xi'}{2} \right\} \quad (\text{by Theorem A.2}) \\ &= 2\Pr \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K \frac{w_k}{N_k} \sum_{n=1}^{N_k} (f(\mathbf{z}'_n^{(k)}) - f(\mathbf{z}_n^{(k)})) \right| > \frac{\xi'}{2} \right\} \\ &= 2\Pr \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K \frac{w_k}{N_k} \sum_{n=1}^{N_k} \epsilon_n^{(k)} (f(\mathbf{z}'_n^{(k)}) - f(\mathbf{z}_n^{(k)})) \right| > \frac{\xi'}{2} \right\} \\ &= 2\Pr \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{f}(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{4} \right\}. \quad (\text{by (63) and (64)}) \end{aligned} \quad (65)$$

Fix a realization of  $\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K$  and let  $\Lambda$  be a  $\xi'/8$ -radius cover of  $\mathcal{F}$  with respect to the  $\ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K)$  norm. Since  $\mathcal{F}$  is composed of the bounded functions with the range  $[a, b]$ , we assume that the same holds for any  $h \in \Lambda$ . If  $f_0$  is the function that achieves the following supremum

$$\sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{f}(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{4},$$

there must be an  $h_0 \in \Lambda$  that satisfies

$$\sum_{k=1}^K \frac{w_k}{2N_k} \left( |f_0(\mathbf{z}_n^{(k)}) - h_0(\mathbf{z}_n^{(k)})| + |f_0(\mathbf{z}_n^{(k)}) - h_0(\mathbf{z}_n^{(k)})| \right) < \frac{\xi'}{8},$$

and meanwhile,

$$\left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{h}_0(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{8}.$$

Therefore, for the realization of  $\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K$ , we arrive at

$$\begin{aligned} & \Pr \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{f}(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{4} \right\} \\ & \leq \Pr \left\{ \sup_{h \in \Lambda} \left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{h}(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{8} \right\}. \end{aligned} \quad (66)$$

Moreover, we denote the event

$$A := \left\{ \Pr \left\{ \sup_{h \in \Lambda} \left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{h}(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{8} \right\} \right\},$$

and let  $\mathbf{1}_A$  be the characteristic function of the event  $A$ . By Fubini's Theorem, we have

$$\begin{aligned} \Pr\{A\} &= \mathbb{E} \left\{ \mathbb{E}_{\vec{\epsilon}} \{ \mathbf{1}_A \} \mid \{\mathbf{Z}_1^{2N_k}\}_{k=1}^K \right\} \\ &= \mathbb{E} \left\{ \Pr \left\{ \sup_{h \in \Lambda} \left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{h}(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{8} \right\} \mid \{\mathbf{Z}_1^{2N_k}\}_{k=1}^K \right\}. \end{aligned} \quad (67)$$

Fix a realization of  $\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K$  again. According to (63), (64) and Theorem A.1, we have

$$\begin{aligned}
& \Pr \left\{ \sup_{h \in \Lambda} \left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{h}(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{8} \right\} \\
& \leq |\Lambda| \max_{h \in \Lambda} \Pr \left\{ \left| \sum_{k=1}^K \frac{w_k}{2N_k} \langle \vec{\epsilon}^{(k)}, \vec{h}(\mathbf{Z}_1^{2N_k}) \rangle \right| > \frac{\xi'}{8} \right\} \\
& = \mathcal{N} \left( \mathcal{F}, \xi'/8, \ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K) \right) \max_{h \in \Lambda} \Pr \left\{ |\mathbf{E}_{\mathbf{w}}^{(S)} h - \mathbf{E}_{\mathbf{w}}^{(S)} h| > \frac{\xi'}{4} \right\} \\
& \leq \mathcal{N} \left( \mathcal{F}, \xi'/8, \ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K) \right) \max_{h \in \Lambda} \Pr \left\{ |\bar{\mathbf{E}}^{(S)} h - \mathbf{E}_{\mathbf{w}}^{(S)} h| + |\bar{\mathbf{E}}^{(S)} h - \mathbf{E}_{\mathbf{w}}^{(S)} h| > \frac{\xi'}{4} \right\} \\
& \leq 2\mathcal{N} \left( \mathcal{F}, \xi'/8, \ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K) \right) \max_{h \in \Lambda} \Pr \left\{ |\bar{\mathbf{E}}^{(S)} h - \mathbf{E}_{\mathbf{w}}^{(S)} h| > \frac{\xi'}{8} \right\} \\
& \leq 4\mathcal{N} \left( \mathcal{F}, \xi'/8, \ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K) \right) \exp \left\{ -\frac{(\prod_{k=1}^K N_k) \left( \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T) \right)^2}{32(b-a)^2 \left( \sum_{k=1}^K w_k^2 (\prod_{i \neq k} N_i) \right)} \right\}, \tag{68}
\end{aligned}$$

where the expectation  $\bar{\mathbf{E}}^{(S)}$  is defined in (57).

The combination of (65), (66) and (68) leads to the result: given an arbitrary  $\xi > D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$  and for any  $\prod_{k=1}^K N_k \geq \frac{8(b-a)^2}{(\xi')^2}$  with  $\xi' = \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ ,

$$\begin{aligned}
& \Pr \left\{ \sup_{f \in \mathcal{F}} |\mathbf{E}^{(T)} f - \mathbf{E}_{\mathbf{w}}^{(S)} f| > \xi \right\} \\
& \leq 8\mathcal{N} \left( \mathcal{F}, \xi'/8, \ell_1^{\mathbf{w}}(\{\mathbf{Z}_1^{2N_k}\}_{k=1}^K) \right) \exp \left\{ -\frac{(\prod_{k=1}^K N_k) \left( \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T) \right)^2}{32(b-a)^2 \left( \sum_{k=1}^K w_k^2 (\prod_{i \neq k} N_i) \right)} \right\} \\
& \leq 8\mathcal{N}_1^{\mathbf{w}} \left( \mathcal{F}, \xi'/8, 2 \sum_{k=1}^K N_k \right) \exp \left\{ -\frac{(\prod_{k=1}^K N_k) \left( \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T) \right)^2}{32(b-a)^2 \left( \sum_{k=1}^K w_k^2 (\prod_{i \neq k} N_i) \right)} \right\}. \tag{69}
\end{aligned}$$

According to (69), letting

$$\epsilon := 8\mathcal{N}_1^{\mathbf{w}} \left( \mathcal{F}, \xi'/8, 2 \sum_{k=1}^K N_k \right) \exp \left\{ -\frac{(\prod_{k=1}^K N_k) \left( \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T) \right)^2}{32(b-a)^2 \left( \sum_{k=1}^K w_k^2 (\prod_{i \neq k} N_i) \right)} \right\},$$

we then arrive at with probability at least  $1 - \epsilon$ ,

$$\sup_{f \in \mathcal{F}} |\mathbf{E}_{\mathbf{w}}^{(S)} f - \mathbf{E}^{(T)} f| \leq D_{\mathcal{F}}^{(\mathbf{w})}(S, T) + \left( \frac{\ln \mathcal{N}_1^{\mathbf{w}}(\mathcal{F}, \xi'/8, 2 \sum_{k=1}^K N_k) - \ln(\epsilon/8)}{\frac{(\prod_{k=1}^K N_k)}{32(b-a)^2 \left( \sum_{k=1}^K w_k^2 (\prod_{i \neq k} N_i) \right)}} \right)^{\frac{1}{2}},$$

where  $\xi' = \xi - D_{\mathcal{F}}^{(\mathbf{w})}(S, T)$ . This completes the proof. ■

## B Proof of Theorem 6.1

Here, we provide the proof of Theorem 6.1. Similar to the situation of domain adaptation with multiple sources, we need to develop the related Hoeffding-type deviation inequality and the symmetrization inequality for domain adaptation combining source and target data.

### B.1 Hoeffding-Type Deviation Inequality

Based on Hoeffding's inequality (Hoeffding, 1963), we derive a deviation inequality for the combination of the source and the target domains.

**Theorem B.1** *Assume that  $\mathcal{F}$  is a function class consisting of bounded functions with the range  $[a, b]$ . Let  $\mathbf{Z}_1^{N_S} := \{\mathbf{z}_n^{(S)}\}_{n=1}^{N_S}$  and  $\bar{\mathbf{Z}}_1^{N_T} := \{\mathbf{z}_n^{(T)}\}_{n=1}^{N_T}$  be sets of i.i.d. samples drawn from the source domain  $\mathcal{Z}^{(S)} \subset \mathbb{R}^L$  and the target domain  $\mathcal{Z}^{(T)} \subset \mathbb{R}^L$ , respectively. For any  $\tau \in [0, 1)$ , define a function  $F_\tau : \mathbb{R}^{L(N_S+N_T)} \rightarrow \mathbb{R}$  as*

$$F_\tau(\mathbf{X}_1^{N_T}, \mathbf{Y}_1^{N_S}) := \tau N_S \sum_{n=1}^{N_T} f(\mathbf{x}_n) + (1 - \tau) N_T \sum_{n=1}^{N_S} f(\mathbf{y}_n), \quad (70)$$

where

$$\mathbf{X}_1^{N_T} := \{\mathbf{x}_1, \dots, \mathbf{x}_{N_T}\} \in (\mathbb{R}^L)^{N_T}; \quad \mathbf{Y}_1^{N_S} := \{\mathbf{y}_1, \dots, \mathbf{y}_{N_S}\} \in (\mathbb{R}^L)^{N_S}.$$

Then, we have for any  $\tau \in [0, 1)$  and any  $\xi > 0$ ,

$$\begin{aligned} & \Pr \left\{ \left| F_\tau(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) - \mathbb{E}^{(*)} F_\tau \right| > \xi \right\} \\ & \leq 2 \exp \left\{ - \frac{2\xi^2}{(b-a)^2 N_S N_T ((1-\tau)^2 N_T + \tau^2 N_S)} \right\}, \end{aligned} \quad (71)$$

where the expectation  $\mathbb{E}^{(*)}$  is taken on both of the source domain  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$ .

In this theorem, we present a deviation inequality for the combination of source and target domains, which is an extension of the classical Hoeffding-type deviation inequality under the assumption of same distribution (see Bousquet *et al.*, 2004, Theorem 1). Compare to the classical result, the resultant deviation inequality (71) allows the random variables to take values from different domains. The two inequalities coincide when the source domain  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$  match, *i.e.*,  $D_{\mathcal{F}}(S, T) = 0$ .

The proof of Theorem B.1 is also processed by a martingale method. Before the formal proof, we introduce some essential notations.

For any  $\tau \in [0, 1)$ , we denote

$$F_S(\mathbf{Z}_1^{N_S}) := (1 - \tau) N_T \sum_{n=1}^{N_S} f(\mathbf{z}_n^{(S)}); \quad F_T(\bar{\mathbf{Z}}_1^{N_T}) := \tau N_S \sum_{n=1}^{N_T} f(\mathbf{z}_n^{(T)}). \quad (72)$$

Recalling (70), it is evident that  $F_\tau(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) = F_S(\mathbf{Z}_1^{N_S}) + F_T(\bar{\mathbf{Z}}_1^{N_T})$ . We then define two random variables:

$$\begin{aligned} S_n &:= \mathbb{E}^{(S)} \{ F_S(\mathbf{Z}_1^{N_S}) | \mathbf{Z}_1^n \}, \quad 0 \leq n \leq N_S; \\ T_n &:= \mathbb{E}^{(T)} \{ F_T(\bar{\mathbf{Z}}_1^{N_T}) | \bar{\mathbf{Z}}_1^n \}, \quad 0 \leq n \leq N_T, \end{aligned} \quad (73)$$

where

$$\begin{aligned} \mathbf{Z}_1^n &= \{\mathbf{z}_1^{(S)}, \dots, \mathbf{z}_n^{(S)}\} \subseteq \mathbf{Z}_1^{N_S} \quad \text{with } \mathbf{Z}_1^0 := \emptyset; \\ \bar{\mathbf{Z}}_1^n &= \{\mathbf{z}_1^{(T)}, \dots, \mathbf{z}_n^{(T)}\} \subseteq \bar{\mathbf{Z}}_1^{N_T} \quad \text{with } \bar{\mathbf{Z}}_1^0 := \emptyset. \end{aligned}$$

It is clear that  $S_0 = \mathbb{E}^{(S)} F_S$ ;  $S_{N_S} = F_S(\mathbf{Z}_1^{N_S})$  and  $T_0 = \mathbb{E}^{(T)} F_T$ ;  $T_{N_T} = F_T(\bar{\mathbf{Z}}_1^{N_T})$ .

According to (70) and (73), we have for any  $1 \leq n \leq N_S$  and any  $\tau \in [0, 1]$ ,

$$\begin{aligned} & S_n - S_{n-1} \\ &= \mathbb{E}^{(S)} \{ F_S(\mathbf{Z}_1^{N_S}) | \mathbf{Z}_1^n \} - \mathbb{E}^{(S)} \{ F_S(\mathbf{Z}_1^{N_S}) | \mathbf{Z}_1^{n-1} \} \\ &= \mathbb{E}^{(S)} \left\{ (1 - \tau) N_T \sum_{n=1}^{N_S} f(\mathbf{z}_n^{(S)}) \middle| \mathbf{Z}_1^n \right\} - \mathbb{E}^{(S)} \left\{ (1 - \tau) N_T \sum_{n=1}^{N_S} f(\mathbf{z}_n^{(S)}) \middle| \mathbf{Z}_1^{n-1} \right\} \\ &= (1 - \tau) N_T \sum_{m=1}^n f(\mathbf{z}_m^{(S)}) + \mathbb{E}^{(S)} \left\{ (1 - \tau) N_T \sum_{m=n+1}^{N_S} f(\mathbf{z}_m^{(S)}) \right\} \\ &\quad - \left( (1 - \tau) N_T \sum_{m=1}^{n-1} f(\mathbf{z}_m^{(S)}) + \mathbb{E}^{(S)} \left\{ (1 - \tau) N_T \sum_{m=n}^{N_S} f(\mathbf{z}_m^{(S)}) \right\} \right) \\ &= (1 - \tau) N_T (f(\mathbf{z}_n^{(S)}) - \mathbb{E}^{(S)} f). \end{aligned} \quad (74)$$

Similarly, we also have for any  $1 \leq n \leq N_T$ ,

$$T_n - T_{n-1} = \tau N_S (f(\mathbf{z}_n^{(T)}) - \mathbb{E}^{(T)} f). \quad (75)$$

We are now ready to prove Theorem B.1.

*Proof of Theorem B.1.* According to (70) and (72), we have

$$\begin{aligned} F_\tau(\mathbf{Z}_1^N) - \mathbb{E}^{(*)} F_\tau &= F_S(\mathbf{Z}_1^{N_S}) + F_T(\bar{\mathbf{Z}}_1^{N_T}) - \mathbb{E}^{(*)} \{ F_S + F_T \} \\ &= F_S(\mathbf{Z}_1^{N_S}) - \mathbb{E}^{(S)} F_S + F_T(\bar{\mathbf{Z}}_1^{N_T}) - \mathbb{E}^{(T)} F_T. \end{aligned} \quad (76)$$

According to Lemma A.1, (74), (75), (76), Markov's inequality, Jensen's inequality and the law of iterated expectation, we have for any  $\alpha > 0$  and any  $\tau \in [0, 1)$ ,

$$\begin{aligned}
& \Pr \{ F_\tau(\mathbf{Z}_1^N) - E^{(*)} F_\tau > \xi \} \\
&= \Pr \left\{ F_S(\mathbf{Z}_1^{N_S}) - E^{(S)} F_S + F_T(\bar{\mathbf{Z}}_1^{N_T}) - E^{(T)} F_T > \xi \right\} \\
&\leq e^{-\alpha \xi} E \left\{ e^{\alpha (F_S(\mathbf{Z}_1^{N_S}) - E^{(S)} F_S + F_T(\bar{\mathbf{Z}}_1^{N_T}) - E^{(T)} F_T)} \right\} \\
&= e^{-\alpha \xi} E \left\{ e^{\alpha (\sum_{n=1}^{N_S} (S_n - S_{n-1}) + \sum_{n=1}^{N_T} (T_n - T_{n-1}))} \right\} \\
&= e^{-\alpha \xi} E \left\{ E \left\{ e^{\alpha (\sum_{n=1}^{N_S} (S_n - S_{n-1}) + \sum_{n=1}^{N_T} (T_n - T_{n-1}))} \middle| \mathbf{Z}_1^{N_S-1} \right\} \right\} \\
&= e^{-\alpha \xi} E \left\{ e^{\alpha (\sum_{n=1}^{N_S-1} (S_n - S_{n-1}) + \sum_{n=1}^{N_T} (T_n - T_{n-1}))} E \left\{ e^{\alpha (S_{N_S} - S_{N_S-1})} \middle| \mathbf{Z}_1^{N_S-1} \right\} \right\} \\
&\leq e^{-\alpha \xi} E \left\{ e^{\alpha (\sum_{n=1}^{N_S-1} (S_n - S_{n-1}) + \sum_{n=1}^{N_T} (T_n - T_{n-1}))} \right\} e^{\frac{(1-\tau)^2 N_T^2 \alpha^2 (b-a)^2}{8}} \\
&= e^{-\alpha \xi} E \left\{ e^{\alpha (\sum_{n=1}^{N_S-1} (S_n - S_{n-1}) + \sum_{n=1}^{N_T-1} (T_n - T_{n-1}))} E \left\{ e^{\alpha (T_{N_T} - T_{N_T-1})} \middle| \bar{\mathbf{Z}}_1^{N_T-1} \right\} \right\} \\
&\quad \times e^{\frac{(1-\tau)^2 N_T^2 \alpha^2 (b-a)^2}{8}} \\
&\leq e^{-\alpha \xi} E \left\{ e^{\alpha (\sum_{n=1}^{N_S-1} (S_n - S_{n-1}) + \sum_{n=1}^{N_T-1} (T_n - T_{n-1}))} \right\} e^{\frac{\tau^2 N_S^2 \alpha^2 (b-a)^2}{8}} e^{\frac{(1-\tau)^2 N_T^2 \alpha^2 (b-a)^2}{8}}. \tag{77}
\end{aligned}$$

Then, we have

$$\Pr \left\{ F_\tau \left( \mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T} \right) - E^{(*)} F_\tau > \xi \right\} \leq e^{\Phi(\alpha) - \alpha \xi}, \tag{78}$$

where

$$\Phi(\alpha) = \frac{\alpha^2 (1-\tau)^2 (b-a)^2 N_S N_T^2}{8} + \frac{\alpha^2 \tau^2 (b-a)^2 N_S^2 N_T}{8}. \tag{79}$$

Similarly, we can arrive at

$$\Pr \left\{ E^{(*)} F_\tau - F_\tau \left( \mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T} \right) > \xi \right\} \leq e^{\Phi(\alpha) - \alpha \xi}. \tag{80}$$

Note that  $\Phi(\alpha) - \alpha \xi$  is a quadratic function with respect to  $\alpha > 0$  and thus the minimum value

$$\min_{\alpha > 0} \{ \Phi(\alpha) - \alpha \xi \}$$

is achieved when

$$\alpha = \frac{4\xi}{(b-a)^2 N_S N_T ((1-\tau)^2 N_T + \tau^2 N_S)}.$$

By combining (78), (79) and (80), we arrive at

$$\Pr \left\{ |F_\tau(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) - E^{(*)} F_\tau| > \xi \right\} \leq 2 \exp \left\{ -\frac{2\xi^2}{(b-a)^2 N_S N_T ((1-\tau)^2 N_T + \tau^2 N_S)} \right\}.$$

This completes the proof. ■



## B.2 Symmetrization Inequality

In the following theorem, we present the symmetrization inequality for domain adaptation combining source and target data.

**Theorem B.2** *Assume that  $\mathcal{F}$  is a function class with the range  $[a, b]$ . Let  $\mathbf{Z}_1^{N_S}$  and  $\mathbf{Z}_1'^{N_S}$  be drawn from the source domain  $\mathcal{Z}^{(S)}$ , and  $\bar{\mathbf{Z}}_1^{N_T}$  and  $\bar{\mathbf{Z}}_1'^{N_T}$  be drawn from the target domain  $\mathcal{Z}^{(T)}$ . Then, for any  $\tau \in [0, 1)$  and given an arbitrary  $\xi > (1 - \tau)D_{\mathcal{F}}(S, T)$ , we have for any  $N_S N_T \geq \frac{8(b-a)^2}{(\xi')^2}$ ,*

$$\Pr \left\{ \sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\tau} f| > \xi \right\} \leq 2 \Pr \left\{ \sup_{f \in \mathcal{F}} |E'_{\tau} f - E_{\tau} f| > \frac{\xi'}{2} \right\} \quad (81)$$

with  $\xi' = \xi - (1 - \tau)D_{\mathcal{F}}(S, T)$ .

This theorem shows that for any  $\xi > (1 - \tau)D_{\mathcal{F}}(S, T)$ , the probability of the event:

$$\sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\tau} f| > \xi$$

can be bounded by using the probability of the event:

$$\sup_{f \in \mathcal{F}} |E'_{\tau} f - E_{\tau} f| > \frac{\xi'}{2}$$

that is only determined by the samples drawn from the source domain  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$ , when  $N_S N_T \geq \frac{8(b-a)^2}{(\xi')^2}$ . Compared to the classical symmetrization result under the assumption of same distribution (see Bousquet *et al.*, 2004), there is a discrepancy term  $(1 - \tau)D_{\mathcal{F}}(S, T)$ . The two results will coincide when the source and the target domains match, *i.e.*,  $D_{\mathcal{F}}(S, T) = 0$ . The following is the proof of Theorem B.2.

*Proof of Theorem B.2.* Let  $\hat{f}$  be the function achieving the supremum:

$$\sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\tau} f|$$

with respect to  $\mathbf{Z}_1^{N_S}$  and  $\bar{\mathbf{Z}}_1^{N_T}$ . According to (9) and (12), we arrive at

$$\begin{aligned} |E^{(T)} \hat{f} - E_{\tau} \hat{f}| &= |\tau E^{(T)} \hat{f} + (1 - \tau) E^{(T)} \hat{f} - (1 - \tau) E^{(S)} \hat{f} + (1 - \tau) E^{(S)} \hat{f} - E_{\tau} \hat{f}| \\ &= |\tau(E^{(T)} \hat{f} - E_{N_T}^{(T)} \hat{f}) + (1 - \tau)(E^{(T)} \hat{f} - E^{(S)} \hat{f}) + (1 - \tau)(E^{(S)} \hat{f} - E_{N_S}^{(S)} \hat{f})| \\ &\leq (1 - \tau)D_{\mathcal{F}}(S, T) + |\tau(E^{(T)} \hat{f} - E_{N_T}^{(T)} \hat{f}) + (1 - \tau)(E^{(S)} \hat{f} - E_{N_S}^{(S)} \hat{f})|, \end{aligned} \quad (82)$$

and thus

$$\begin{aligned} &\Pr \left\{ |E^{(T)} \hat{f} - E_{\tau} \hat{f}| > \xi \right\} \\ &\leq \Pr \left\{ (1 - \tau)D_{\mathcal{F}}(S, T) + |\tau(E^{(T)} \hat{f} - E_{N_T}^{(T)} \hat{f}) + (1 - \tau)(E^{(S)} \hat{f} - E_{N_S}^{(S)} \hat{f})| > \xi \right\}, \end{aligned} \quad (83)$$

where

$$\mathbf{E}_{N_T}^{(T)} \hat{f} := \frac{1}{N_T} \sum_{n=1}^{N_T} \hat{f}(\mathbf{z}_n^{(T)}); \quad \mathbf{E}_{N_S}^{(S)} \hat{f} := \frac{1}{N_S} \sum_{n=1}^{N_S} \hat{f}(\mathbf{z}_n^{(S)}). \quad (84)$$

Let

$$\xi' = \xi - (1 - \tau)D_{\mathcal{F}}(S, T) \quad (85)$$

and denote  $\wedge$  as the conjunction of two events. According to the triangle inequality, we have

$$\begin{aligned} & \left( \mathbf{1}_{\left\{ |\tau(\mathbf{E}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f})| > \xi' \right\}} \right) \left( \mathbf{1}_{\left\{ |\tau(\mathbf{E}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f})| < \frac{\xi'}{2} \right\}} \right) \\ &= \mathbf{1}_{\left\{ |\tau(\mathbf{E}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f})| > \xi' \right\}} \wedge \left\{ |\tau(\mathbf{E}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f})| < \frac{\xi'}{2} \right\} \\ &\leq \mathbf{1}_{\left\{ |\tau(\mathbf{E}_{N_T}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}_{N_S}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f})| > \frac{\xi'}{2} \right\}}. \end{aligned}$$

Then, taking the expectation with respect to  $\mathbf{Z}_1^{N_S}$  and  $\overline{\mathbf{Z}}_1^{N_T}$  gives

$$\begin{aligned} & \left( \mathbf{1}_{\left\{ |\tau(\mathbf{E}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f})| > \xi' \right\}} \right) \\ & \times \Pr' \left\{ \left| \tau(\mathbf{E}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f}) \right| < \frac{\xi'}{2} \right\} \\ & \leq \Pr' \left\{ \left| \tau(\mathbf{E}_{N_T}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}_{N_S}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f}) \right| > \frac{\xi'}{2} \right\}. \end{aligned} \quad (86)$$

By Chebyshev's inequality, since  $\mathbf{Z}_1^{N_S} = \{\mathbf{z}_n^{(S)}\}_{n=1}^{N_S}$  and  $\overline{\mathbf{Z}}_1^{N_T} = \{\mathbf{z}_n^{(T)}\}_{n=1}^{N_T}$  are sets of i.i.d. samples drawn from the source domain  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$  respectively, we have for any  $\xi' > 0$  and any  $\tau \in [0, 1]$ ,

$$\begin{aligned} & \Pr' \left\{ \left| \tau(\mathbf{E}^{(T)} \hat{f} - \mathbf{E}_{N_T}^{(T)} \hat{f}) + (1-\tau)(\mathbf{E}^{(S)} \hat{f} - \mathbf{E}_{N_S}^{(S)} \hat{f}) \right| \geq \frac{\xi'}{2} \right\} \\ & \leq \Pr' \left\{ \frac{\tau}{N_T} \sum_{n=1}^{N_T} |\mathbf{E}^{(T)} \hat{f} - \hat{f}(\mathbf{z}_n^{(T)})| + \frac{1-\tau}{N_S} \sum_{n=1}^{N_S} |\mathbf{E}^{(S)} \hat{f} - \hat{f}(\mathbf{z}_n^{(S)})| \geq \frac{\xi'}{2} \right\} \\ & \leq \frac{4\mathbf{E} \left\{ \tau N_S N_T (\mathbf{E}^{(T)} \hat{f} - \hat{f}(\mathbf{z}'^{(T)}))^2 + (1-\tau) N_S N_T (\mathbf{E}^{(S)} \hat{f} - \hat{f}(\mathbf{z}'^{(S)}))^2 \right\}}{N_S^2 N_T^2 (\xi')^2} \\ & \leq \frac{4\mathbf{E} \left\{ \tau N_S N_T (b-a)^2 + (1-\tau) N_S N_T (b-a)^2 \right\}}{N_S^2 N_T^2 (\xi')^2} \\ & = \frac{4(b-a)^2}{N_S N_T (\xi')^2}, \end{aligned} \quad (87)$$

where  $\mathbf{z}'^{(S)}$  and  $\mathbf{z}'^{(T)}$  stand for the ghost random variables taking values from the source domain  $\mathcal{Z}^{(S)}$  and the target domain  $\mathcal{Z}^{(T)}$ , respectively.

Subsequently, according to (86) and (87), we have for any  $\xi' > 0$ ,

$$\begin{aligned} & \Pr' \left\{ \left| \tau(E_{N_T}^{(T)} \hat{f} - E_{N_T}^{(T)} \hat{f}) + (1 - \tau)(E_{N_S}^{(S)} \hat{f} - E_{N_S}^{(S)} \hat{f}) \right| > \frac{\xi'}{2} \right\} \\ & \geq \left( \mathbf{1}_{\left\{ \left| \tau(E^{(T)} \hat{f} - E_{N_T}^{(T)} \hat{f}) + (1 - \tau)(E^{(S)} \hat{f} - E_{N_S}^{(S)} \hat{f}) \right| > \xi' \right\}} \right) \left( 1 - \frac{4(b - a)^2}{N_S N_T (\xi')^2} \right). \end{aligned} \quad (88)$$

According to (83), (85) and (88), by letting

$$\frac{4(b - a)^2}{N_S N_T (\xi')^2} \leq \frac{1}{2},$$

and taking the expectation with respect to  $\mathbf{Z}_1^{N_S}$  and  $\overline{\mathbf{Z}}_1^{N_T}$ , we have for any  $\xi' > 0$ ,

$$\begin{aligned} & \Pr \left\{ |E^{(T)} \hat{f} - E_{\tau} \hat{f}| > \xi \right\} \\ & \leq \Pr \left\{ \left| \tau(E^{(T)} \hat{f} - E_{N_T}^{(T)} \hat{f}) + (1 - \tau)(E^{(S)} \hat{f} - E_{N_S}^{(S)} \hat{f}) \right| > \xi' \right\} \\ & \leq 2 \Pr \left\{ \left| \tau(E_{N_T}^{(T)} \hat{f} - E_{N_T}^{(T)} \hat{f}) + (1 - \tau)(E_{N_S}^{(S)} \hat{f} - E_{N_S}^{(S)} \hat{f}) \right| > \frac{\xi'}{2} \right\} \\ & = 2 \Pr \left\{ |E'_{\tau} \hat{f} - E_{\tau} \hat{f}| > \frac{\xi'}{2} \right\} \end{aligned} \quad (89)$$

with  $\xi' = \xi - (1 - \tau)D_{\mathcal{F}}(S, T)$ . This completes the proof.  $\blacksquare$

We are now ready to prove Theorem 6.1.

### B.3 Proof of Theorem 6.1

*Proof of Theorem 6.1.* Consider  $\{\epsilon_n\}_{n=1}^N$  as independent Rademacher random variables, *i.e.*, independent  $\{\pm 1\}$ -valued random variables with equal probability of taking either value. Given  $\{\epsilon_n\}_{n=1}^{N_S}$ ,  $\{\epsilon_n\}_{n=1}^{N_T}$ ,  $\mathbf{Z}_1^{2N_S}$  and  $\overline{\mathbf{Z}}_1^{2N_T}$ , denote

$$\begin{aligned} \vec{\epsilon}_S & := (\epsilon_1, \dots, \epsilon_{N_S}, -\epsilon_1, \dots, -\epsilon_{N_S}) \in \{\pm 1\}^{2N_S}; \\ \vec{\epsilon}_T & := (\epsilon_1, \dots, \epsilon_{N_T}, -\epsilon_1, \dots, -\epsilon_{N_T}) \in \{\pm 1\}^{2N_T}, \end{aligned} \quad (90)$$

and for any  $f \in \mathcal{F}$ ,

$$\begin{aligned} \vec{f}(\mathbf{Z}_1^{2N_S}) & := (f(\mathbf{z}'_1), \dots, f(\mathbf{z}'_{N_S}), f(\mathbf{z}_1), \dots, f(\mathbf{z}_{N_S})) \in [a, b]^{2N_S}; \\ \vec{f}(\mathbf{Z}_1^{2N_T}) & := (f(\mathbf{z}'_1), \dots, f(\mathbf{z}'_{N_T}), f(\mathbf{z}_1), \dots, f(\mathbf{z}_{N_T})) \in [a, b]^{2N_T}. \end{aligned} \quad (91)$$

We also denote

$$\begin{aligned} \mathbf{Z} & := \left\{ \overline{\mathbf{Z}}_1^{2N_T}, \mathbf{Z}_1^{2N_S} \right\} \in (\mathcal{Z}^{(T)})^{2N_T} \times (\mathcal{Z}^{(S)})^{2N_S}; \\ \vec{\epsilon} & := \left( \underbrace{\vec{\epsilon}_T, \dots, \vec{\epsilon}_T}_{N_S}, \underbrace{\vec{\epsilon}_S, \dots, \vec{\epsilon}_S}_{N_T} \right) \in \{\pm 1\}^{4N_S N_T}; \\ \vec{f}(\mathbf{Z}) & := \left( \underbrace{\vec{f}(\overline{\mathbf{Z}}_1^{2N_T}), \dots, \vec{f}(\overline{\mathbf{Z}}_1^{2N_T})}_{N_S}, \underbrace{\vec{f}(\mathbf{Z}_1^{2N_S}), \dots, \vec{f}(\mathbf{Z}_1^{2N_S})}_{N_T} \right) \in [a, b]^{4N_S N_T}. \end{aligned} \quad (92)$$

According to (6), (85) and Theorem B.2, for any  $\tau \in [0, 1)$  and given an arbitrary  $\xi > (1 - \tau)D_{\mathcal{F}}(S, T)$ , we have for any  $N_S N_T \geq \frac{8(b-a)^2}{\xi'^2}$  with  $\xi' = \xi - (1 - \tau)D_{\mathcal{F}}(S, T)$ ,

$$\begin{aligned}
& \Pr \left\{ \sup_{f \in \mathcal{F}} |E^{(T)} f - E_{\tau} f| > \xi \right\} \\
& \leq 2 \Pr \left\{ \sup_{f \in \mathcal{F}} |E'_{\tau} f - E_{\tau} f| > \frac{\xi'}{2} \right\} \quad (\text{by Theorem B.2}) \\
& = 2 \Pr \left\{ \sup_{f \in \mathcal{F}} \left| \frac{\tau}{N_T} \sum_{n=1}^{N_T} (f(\mathbf{z}'^{(T)}_n) - f(\mathbf{z}^{(T)}_n)) + \frac{1-\tau}{N_S} \sum_{n=1}^{N_S} (f(\mathbf{z}'^{(S)}_n) - f(\mathbf{z}^{(S)}_n)) \right| > \frac{\xi'}{2} \right\} \\
& = 2 \Pr \left\{ \sup_{f \in \mathcal{F}} \left| \frac{\tau}{N_T} \sum_{n=1}^{N_T} \epsilon_n (f(\mathbf{z}'^{(T)}_n) - f(\mathbf{z}^{(T)}_n)) + \frac{1-\tau}{N_S} \sum_{n=1}^{N_S} \epsilon_n (f(\mathbf{z}'^{(S)}_n) - f(\mathbf{z}^{(S)}_n)) \right| > \frac{\xi'}{2} \right\} \\
& = 2 \Pr \left\{ \sup_{f \in \mathcal{F}} \left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{f}(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{f}(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{4} \right\}. \tag{93}
\end{aligned}$$

Given a  $\tau \in [0, 1)$ , fix a realization of  $\mathbf{Z}$  and let  $\Lambda$  be a  $\xi'/8$ -radius cover of  $\mathcal{F}$  with respect to the  $\ell_1^{\tau}(\mathbf{Z})$  norm. Since  $\mathcal{F}$  is composed of bounded functions with the range  $[a, b]$ , we assume that the same holds for any  $h \in \Lambda$ . If  $f_0$  is the function that achieves the following supremum

$$\sup_{f \in \mathcal{F}} \left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{f}(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{f}(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{4},$$

there must be an  $h_0 \in \Lambda$  that satisfies that

$$\begin{aligned}
& \frac{\tau}{2N_T} \sum_{n=1}^{N_T} (|f_0(\mathbf{z}'^{(T)}_n) - h_0(\mathbf{z}'^{(T)}_n)| + |f_0(\mathbf{z}^{(T)}_n) - h_0(\mathbf{z}^{(T)}_n)|) \\
& + \frac{1-\tau}{2N_S} \sum_{n=1}^{N_S} (|f_0(\mathbf{z}'^{(S)}_n) - h_0(\mathbf{z}'^{(S)}_n)| + |f_0(\mathbf{z}^{(S)}_n) - h_0(\mathbf{z}^{(S)}_n)|) < \frac{\xi'}{8},
\end{aligned}$$

and meanwhile,

$$\left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{h}_0(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{h}_0(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{8}.$$

Therefore, for the realization of  $\mathbf{Z}$ , we arrive at

$$\begin{aligned}
& \Pr \left\{ \sup_{f \in \mathcal{F}} \left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{f}(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{f}(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{4} \right\} \\
& \leq \Pr \left\{ \sup_{h \in \Lambda} \left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{h}(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{h}(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{8} \right\}. \tag{94}
\end{aligned}$$

Moreover, we denote the event

$$A := \left\{ \Pr \left\{ \sup_{h \in \Lambda} \left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{h}(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{h}(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{8} \right\} \right\},$$

and let  $\mathbf{1}_A$  be the characteristic function of the event  $A$ . By Fubini's Theorem, we have

$$\begin{aligned} \Pr\{A\} &= \mathbb{E}\left\{\mathbb{E}_{\vec{\epsilon}}\{\mathbf{1}_A\} \mid \mathbf{Z}\right\} \\ &= \mathbb{E}\left\{\Pr\left\{\sup_{h \in \Lambda} \left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{h}(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{h}(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{8} \right\} \mid \mathbf{Z}\right\}. \end{aligned} \quad (95)$$

Fix a realization of  $\mathbf{Z}$  again. According to (21), (90), (91) and Theorem B.1, for any  $\tau \in [0, 1]$  and given an arbitrary  $\xi' > 0$ , we have for any  $N_S N_T \geq \frac{8(b-a)^2}{(\xi')^2}$ ,

$$\begin{aligned} &\Pr\left\{\sup_{h \in \Lambda} \left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{h}(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{h}(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{8} \right\} \\ &\leq |\Lambda| \max_{h \in \Lambda} \Pr\left\{\left| \frac{\tau}{2N_T} \langle \vec{\epsilon}_T, \vec{h}(\bar{\mathbf{Z}}_1^{2N_T}) \rangle + \frac{1-\tau}{2N_S} \langle \vec{\epsilon}_S, \vec{h}(\mathbf{Z}_1^{2N_S}) \rangle \right| > \frac{\xi'}{8} \right\} \\ &= \mathcal{N}(\mathcal{F}, \xi'/8, \ell_1^T(\mathbf{Z})) \max_{h \in \Lambda} \Pr\left\{|\mathbb{E}'_{\tau} h - \mathbb{E}_{\tau} h| > \frac{\xi'}{4} \right\} \\ &\leq \mathcal{N}(\mathcal{F}, \xi'/8, \ell_1^T(\mathbf{Z})) \max_{h \in \Lambda} \Pr\left\{|\tilde{\mathbb{E}} h - \mathbb{E}'_{\tau} h| + |\tilde{\mathbb{E}} h - \mathbb{E}_{\tau} h| > \frac{\xi'}{4} \right\} \\ &\leq 2\mathcal{N}(\mathcal{F}, \xi'/8, \ell_1^T(\mathbf{Z})) \max_{h \in \Lambda} \Pr\left\{|\tilde{\mathbb{E}} h - \mathbb{E}_{\tau} h| > \frac{\xi'}{8} \right\} \\ &\leq 4\mathcal{N}(\mathcal{F}, \xi'/8, \ell_1^T(\mathbf{Z})) \exp\left\{-\frac{N_S N_T (\xi - (1-\tau)D_{\mathcal{F}}(S, T))^2}{32(b-a)^2 ((1-\tau)^2 N_T + \tau^2 N_S)}\right\}, \end{aligned} \quad (96)$$

where  $\tilde{\mathbb{E}} h := \tau \mathbb{E}^{(T)} h + (1-\tau) \mathbb{E}^{(S)} h$ .

The combination of (22), (93), (94) and (96) leads to the following result: for any  $\tau \in [0, 1]$  and given an arbitrary  $\xi > (1-\tau)D_{\mathcal{F}}(S, T)$ , we have for any  $N_S N_T \geq \frac{8(b-a)^2}{(\xi')^2}$ ,

$$\begin{aligned} &\Pr\left\{\sup_{f \in \mathcal{F}} |\mathbb{E}^{(T)} f - \mathbb{E}_{\tau} f| > \xi \right\} \\ &\leq 8\mathcal{N}(\mathcal{F}, \xi'/8, \ell_1^T(\mathbf{Z})) \exp\left\{-\frac{N_S N_T (\xi - (1-\tau)D_{\mathcal{F}}(S, T))^2}{32(b-a)^2 ((1-\tau)^2 N_T + \tau^2 N_S)}\right\} \\ &\leq 8\mathcal{N}_1^T(\mathcal{F}, \xi'/8, 2(N_S + N_T)) \exp\left\{-\frac{N_S N_T (\xi - (1-\tau)D_{\mathcal{F}}(S, T))^2}{32(b-a)^2 ((1-\tau)^2 N_T + \tau^2 N_S)}\right\}. \end{aligned} \quad (97)$$

According to (97), letting

$$\epsilon := 8\mathcal{N}_1^T(\mathcal{F}, \xi'/8, 2(N_S + N_T)) \exp\left\{-\frac{N_S N_T (\xi - (1-\tau)D_{\mathcal{F}}(S, T))^2}{32(b-a)^2 ((1-\tau)^2 N_T + \tau^2 N_S)}\right\},$$

we have given an arbitrary  $\xi > (1-\tau)D_{\mathcal{F}}(S, T)$  and for any  $N_S N_T \geq \frac{8(b-a)^2}{(\xi')^2}$ , with probability at least  $1 - \epsilon$ ,

$$\sup_{f \in \mathcal{F}} |\mathbb{E}_{\tau} f - \mathbb{E}^{(T)} f| \leq (1-\tau)D_{\mathcal{F}}(S, T) + \left( \frac{\ln \mathcal{N}_1^T(\mathcal{F}, \xi'/8, 2(N_S + N_T)) - \ln(\epsilon/8)}{\frac{N_S N_T}{32(b-a)^2 ((1-\tau)^2 N_T + \tau^2 N_S)}} \right)^{\frac{1}{2}}.$$

This completes the proof. ■

## C Proofs of Theorems 5.2 & 6.2

In this appendix, we will prove Theorem 5.2 and Theorem 6.2. In order to achieve the proofs, we need to generalize the classical McDiarmid's inequality (see Bousquet *et al.*, 2004, Theorem 6) to a more general setting where independent random variables can independently take values from different domains.

### C.1 Generalized McDiarmid's inequality

The following is the classical McDiarmid's inequality that is one of the most frequently used deviation inequalities in statistical learning theory and has been widely used to obtain generalization bounds based on the Rademacher complexity under the assumption of same distribution (see Bousquet *et al.*, 2004, Theorem 6).

**Theorem C.1 (McDiarmid's Inequality)** *Let  $\mathbf{z}_1, \dots, \mathbf{z}_N$  be  $N$  independent random variables taking value from the domain  $\mathcal{Z}$ . Assume that the function  $H : \mathcal{Z} \rightarrow \mathbb{R}$  satisfies the condition of bounded difference: for all  $1 \leq n \leq N$ ,*

$$\sup_{\mathbf{z}_1, \dots, \mathbf{z}_N, \mathbf{z}'_n} \left| H(\mathbf{z}_1, \dots, \mathbf{z}_n, \dots, \mathbf{z}_N) - H(\mathbf{z}_1, \dots, \mathbf{z}'_n, \dots, \mathbf{z}_N) \right| \leq c_n. \quad (98)$$

Then, for any  $\xi > 0$

$$\Pr \left\{ H(\mathbf{z}_1, \dots, \mathbf{z}_n, \dots, \mathbf{z}_N) - \mathbb{E} \left\{ H(\mathbf{z}_1, \dots, \mathbf{z}_n, \dots, \mathbf{z}_N) \right\} \geq \xi \right\} \leq \exp \left\{ -2\xi^2 / \sum_{n=1}^N c_n^2 \right\}.$$

As shown in Theorem C.1, the classical McDiarmid's inequality is valid under the condition that random variables  $\mathbf{z}_1, \dots, \mathbf{z}_N$  are independent and drawn from the same domain. Next, we generalize this inequality to a more general setting, where the independent random variables can take values from different domains.

**Theorem C.2** *Given independent domains  $\mathcal{Z}^{(S_k)}$  ( $1 \leq k \leq K$ ), for any  $1 \leq k \leq K$ , let  $\mathbf{Z}_1^{N_k} := \{\mathbf{z}_n^{(S_k)}\}_{n=1}^{N_k}$  be  $N_k$  independent random variables taking values from the domain  $\mathcal{Z}^{(S_k)}$ . Assume that the function  $H : (\mathcal{Z}^{(S_1)})^{N_1} \times \dots \times (\mathcal{Z}^{(S_K)})^{N_K} \rightarrow \mathbb{R}$  satisfies the condition of bounded difference: for all  $1 \leq k \leq K$  and  $1 \leq n \leq N_k$ ,*

$$\sup_{\mathbf{z}_1^{N_1}, \dots, \mathbf{z}_1^{N_K}, \mathbf{z}'_n^{(S_k)}} \left| H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{z}_1^{(S_k)}, \dots, \mathbf{z}_n^{(S_k)}, \dots, \mathbf{z}_{N_k}^{(S_k)}, \mathbf{Z}_1^{N_{k+1}}, \dots, \mathbf{Z}_1^{N_K}) \right. \\ \left. - H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{z}_1^{(S_k)}, \dots, \mathbf{z}'_n^{(S_k)}, \dots, \mathbf{z}_{N_k}^{(S_k)}, \mathbf{Z}_1^{N_{k+1}}, \dots, \mathbf{Z}_1^{N_K}) \right| \leq c_n^{(k)}. \quad (99)$$

Then, for any  $\xi > 0$

$$\Pr \left\{ H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K}) - \mathbb{E} \left\{ H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K}) \right\} \geq \xi \right\} \leq \exp \left\{ -2\xi^2 / \sum_{k=1}^K \sum_{n=1}^{N_k} (c_n^{(k)})^2 \right\}.$$

**Proof.** Define a random variable

$$T_n^{(k)} := \mathbb{E} \left\{ H(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) | \mathbf{Z}_1^{N_1}, \mathbf{Z}_1^{N_2}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{Z}_1^n \right\}, \quad 1 \leq k \leq K, \quad 0 \leq n \leq N_k, \quad (100)$$

where

$$\mathbf{Z}_1^n = \{\mathbf{z}_1^{(k)}, \mathbf{z}_2^{(k)}, \dots, \mathbf{z}_n^{(k)}\} \subseteq \mathbf{Z}_1^{N_k}, \text{ and } \mathbf{Z}_1^0 = \emptyset.$$

It is clear that

$$T_0^{(1)} = \mathbb{E}\{H(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K)\} \text{ and } T_{N_K}^{(K)} = H(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K),$$

and thus

$$H(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K) - \mathbb{E}\{H(\{\mathbf{Z}_1^{N_k}\}_{k=1}^K)\} = T_{N_K}^{(K)} - T_0^{(1)} = \sum_{k=1}^K \sum_{n=1}^{N_k} (T_n^{(k)} - T_{n-1}^{(k)}). \quad (101)$$

Denote for any  $1 \leq k \leq K$  and  $1 \leq n \leq N_k$ ,

$$U_n^{(k)} = \sup_{\mu} \left\{ T_n^{(k)} \big|_{\mathbf{z}_n^{(k)} = \mu} - T_{n-1}^{(k)} \right\};$$

$$L_n^{(k)} = \inf_{\nu} \left\{ T_n^{(k)} \big|_{\mathbf{z}_n^{(k)} = \nu} - T_{n-1}^{(k)} \right\}.$$

It follows from the definition of (100) that  $L_n^{(k)} \leq (T_n^{(k)} - T_{n-1}^{(k)}) \leq U_n^{(k)}$  and thus results in

$$T_n^{(k)} - T_{n-1}^{(k)} \leq U_n^{(k)} - L_n^{(k)} = \sup_{\mu, \nu} \left\{ T_n^{(k)} \big|_{\mathbf{z}_n^{(k)} = \mu} - T_n^{(k)} \big|_{\mathbf{z}_n^{(k)} = \nu} \right\} \leq c_n^{(k)}. \quad (102)$$

Moreover, by the law of iterated expectation, we also have for any  $1 \leq k \leq K$  and  $1 \leq n \leq N_k$

$$\mathbb{E} \left\{ T_n^{(k)} - T_{n-1}^{(k)} | \mathbf{Z}_1^{N_1}, \mathbf{Z}_1^{N_2}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{Z}_1^{n-1} \right\} = 0. \quad (103)$$

According to Hoeffding inequality (see Hoeffding, 1963), given an  $\alpha > 0$ , the condition (99) leads to for any  $1 \leq k \leq K$  and  $1 \leq n \leq N_k$ ,

$$\mathbb{E} \left\{ e^{\alpha(T_n^{(k)} - T_{n-1}^{(k)})} | \mathbf{Z}_1^{N_1}, \mathbf{Z}_1^{N_2}, \dots, \mathbf{Z}_1^{N_{k-1}}, \mathbf{Z}_1^{n-1} \right\} \leq e^{\alpha^2 (c_n^{(k)})^2 / 8}. \quad (104)$$

Subsequently, according to Markov's inequality, (101), (102), (103) and (104), we have for any  $\alpha > 0$ ,

$$\begin{aligned}
& \Pr \left\{ H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K}) - \mathbb{E} \left\{ H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K}) \right\} \geq \xi \right\} \\
& \leq e^{-\alpha \xi} \mathbb{E} \left\{ e^{\alpha \left( H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K}) - \mathbb{E} \left\{ H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K}) \right\} \right)} \right\} \\
& = e^{-\alpha \xi} \mathbb{E} \left\{ e^{\alpha \left( \sum_{k=1}^K \sum_{n=1}^{N_k} (T_n^{(k)} - T_{n-1}^{(k)}) \right)} \right\} \\
& = e^{-\alpha \xi} \mathbb{E} \left\{ \mathbb{E} \left\{ e^{\alpha \left( \sum_{k=1}^K \sum_{n=1}^{N_k} (T_n^{(k)} - T_{n-1}^{(k)}) \right)} \middle| \mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_{K-1}}, \mathbf{Z}_1^{N_K-1} \right\} \right\} \\
& = e^{-\alpha \xi} \mathbb{E} \left\{ e^{\alpha \left( \sum_{k=1}^{K-1} \sum_{n=1}^{N_k} (T_n^{(k)} - T_{n-1}^{(k)}) + \sum_{n=1}^{N_K-1} (T_n^{(K)} - T_{n-1}^{(K)}) \right)} \left\{ e^{\alpha (T_{N_K}^{(K)} - T_{N_K-1}^{(K)})} \middle| \mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_{K-1}}, \mathbf{Z}_1^{N_K-1} \right\} \right\} \\
& \leq e^{-\alpha \xi} \mathbb{E} \left\{ e^{\alpha \left( \sum_{k=1}^{K-1} \sum_{n=1}^{N_k} (T_n^{(k)} - T_{n-1}^{(k)}) + \sum_{n=1}^{N_K-1} (T_n^{(K)} - T_{n-1}^{(K)}) \right)} \right\} e^{\alpha^2 (c_{N_K}^{(K)})^2 / 8} \\
& \leq e^{-\alpha \xi} \prod_{k=1}^K \prod_{n=1}^{N_k} \exp \left\{ \frac{\alpha^2 (c_n^{(k)})^2}{8} \right\} \\
& = \exp \left\{ -\alpha \xi + \alpha^2 \sum_{k=1}^K \sum_{n=1}^{N_k} \frac{(c_n^{(k)})^2}{8} \right\}.
\end{aligned}$$

The above bound is minimized by setting

$$\alpha^* = \frac{4\xi}{\sum_{k=1}^K \sum_{n=1}^{N_k} (c_n^{(k)})^2},$$

and its minimum value is

$$\exp \left\{ -2\xi^2 / \sum_{k=1}^K \sum_{n=1}^{N_k} (c_n^{(k)})^2 \right\}.$$

This completes the proof. ■

## C.2 Proof of Theorem 5.2

By using Theorem C.2, we prove Theorem 5.2 as follows:

**Proof of Theorem 5.2** Assume that  $\mathcal{F}$  is a function class  $\mathcal{F}$  consisting of bounded functions with the range  $[a, b]$ . Let sample sets  $\{\mathbf{Z}_1^{N_k}\}_{k=1}^K := \{\{\mathbf{z}_n^{(k)}\}_{n=1}^{N_k}\}_{k=1}^K$  be drawn from multiple sources  $\mathcal{Z}^{(S_k)}$  ( $1 \leq k \leq K$ ), respectively. Given a choice of  $\mathbf{w} \in [0, 1]^K$  with  $\sum_{k=1}^K w_k = 1$ , denote

$$H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K}) := \sup_{f \in \mathcal{F}} |\mathbb{E}_{\mathbf{w}}^{(S)} f - \mathbb{E}^{(T)} f|. \quad (105)$$

By (1), we have

$$H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K}) = \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K w_k (\mathbb{E}_{N_k}^{(S)} f - \mathbb{E}^{(T)} f) \right|, \quad (106)$$

where  $\mathbb{E}_{N_k}^{(S)} f = \frac{1}{N_k} \sum_{n=1}^{N_k} f(\mathbf{z}_n^{(k)})$ . Therefore, it is clear that such  $H(\mathbf{Z}_1^{N_1}, \dots, \mathbf{Z}_1^{N_K})$  satisfies the condition of bounded difference with

$$c_n^{(k)} = \frac{(b-a)w_k}{N_k}$$



for all  $1 \leq k \leq K$  and  $1 \leq n \leq N_k$ . Thus, according to Theorem C.2, we have for any  $\xi > 0$ ,

$$\Pr\left\{H(\mathbf{Z}_1^{N_k}, \dots, \mathbf{Z}_1^{N_k}) - \mathbb{E}^{(S)}\{H(\mathbf{Z}_1^{N_k}, \dots, \mathbf{Z}_1^{N_k})\} \geq \xi\right\} \leq \exp\left\{\frac{-2\xi^2}{\sum_{k=1}^K \frac{(b-a)^2 w_k^2}{N_k}}\right\},$$

which can be equivalently rewritten as with probability at least  $1 - \epsilon$ ,

$$\begin{aligned} H(\mathbf{Z}_1^{N_k}, \dots, \mathbf{Z}_1^{N_k}) &\leq \mathbb{E}^{(S)}\{H(\mathbf{Z}_1^{N_k}, \dots, \mathbf{Z}_1^{N_k})\} + \sqrt{\sum_{k=1}^K \frac{(b-a)^2 w_k^2 \ln(1/\epsilon)}{2N_k}} \\ &= \mathbb{E}^{(S)}\left\{\sup_{f \in \mathcal{F}} \left|\sum_{k=1}^K w_k (\mathbb{E}_{N_k}^{(S)} f - \mathbb{E}^{(T)} f)\right|\right\} + \sqrt{\sum_{k=1}^K \frac{(b-a)^2 w_k^2 \ln(1/\epsilon)}{2N_k}} \\ &= \mathbb{E}^{(S)}\left\{\sup_{f \in \mathcal{F}} \left|\sum_{k=1}^K w_k (\mathbb{E}_{N_k}^{(S)} f - \mathbb{E}^{(S_k)} f + \mathbb{E}^{(S_k)} f - \mathbb{E}^{(T)} f)\right|\right\} \\ &\quad + \sqrt{\sum_{k=1}^K \frac{(b-a)^2 w_k^2 \ln(1/\epsilon)}{2N_k}} \\ &\leq \mathbb{E}^{(S)}\left\{\sup_{f \in \mathcal{F}} \left|\sum_{k=1}^K w_k (\mathbb{E}_{N_k}^{(S)} f - \mathbb{E}^{(S_k)} f)\right|\right\} + \sum_{k=1}^K w_k \sup_{f \in \mathcal{F}} |\mathbb{E}^{(S_k)} f - \mathbb{E}^{(T)} f| \\ &\quad + \sqrt{\sum_{k=1}^K \frac{(b-a)^2 w_k^2 \ln(1/\epsilon)}{2N_k}} \\ &\leq \mathbb{E}^{(S)}\left\{\sup_{f \in \mathcal{F}} \left|\sum_{k=1}^K w_k (\mathbb{E}_{N_k}^{(S)} f - \mathbb{E}^{(S_k)} f)\right|\right\} + D_{\mathcal{F}}^{(\mathbf{w})}(S, T) \quad [\text{see (26)}] \\ &\quad + \sqrt{\sum_{k=1}^K \frac{(b-a)^2 w_k^2 \ln(1/\epsilon)}{2N_k}}. \end{aligned} \tag{107}$$

Next, according to (23) and (106), we have

$$\begin{aligned}
& \mathbb{E}^{(S)} \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K w_k (\mathbb{E}_{N_k}^{(S)} f - \mathbb{E}^{(S_k)} f) \right| \\
&= \mathbb{E}^{(S)} \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K w_k (\mathbb{E}_{N_k}^{(S)} f - \mathbb{E}'^{(S_k)} \{\mathbb{E}_{N_k}^{(S)} f\}) \right| \\
&\leq \mathbb{E}^{(S)} \mathbb{E}'^{(S)} \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K w_k (\mathbb{E}_{N_k}^{(S)} f - \mathbb{E}_{N_k}^{(S)} f) \right| \\
&= \mathbb{E}^{(S)} \mathbb{E}'^{(S)} \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K w_k \frac{1}{N_k} \sum_{n=1}^{N_k} (f(\mathbf{z}_n^{(k)}) - f(\mathbf{z}'_n^{(k)})) \right| \\
&\leq \mathbb{E}^{(S)} \mathbb{E}'^{(S)} \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K w_k \frac{1}{N_k} \sum_{n=1}^{N_k} \sigma_n^{(k)} (f(\mathbf{z}_n^{(k)}) - f(\mathbf{z}'_n^{(k)})) \right| \\
&\leq 2 \mathbb{E}^{(S)} \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \left| \sum_{k=1}^K w_k \frac{1}{N_k} \sum_{n=1}^{N_k} \sigma_n^{(k)} f(\mathbf{z}_n^{(k)}) \right| \\
&\leq 2 \mathbb{E}^{(S)} \mathbb{E}_\sigma \sum_{k=1}^K w_k \frac{1}{N_k} \sup_{f \in \mathcal{F}} \left| \sum_{n=1}^{N_k} \sigma_n^{(k)} f(\mathbf{z}_n^{(k)}) \right| \\
&= 2 \sum_{k=1}^K w_k \mathcal{R}^{(k)}(\mathcal{F}). \tag{108}
\end{aligned}$$

By combining (105), (107) and (108), we obtain with probability at least  $1 - \epsilon$

$$\sup_{f \in \mathcal{F}} |\mathbb{E}_{\mathbf{w}}^{(S)} f - \mathbb{E}^{(T)} f| \leq D_{\mathcal{F}}^{(\mathbf{w})}(S, T) + 2 \sum_{k=1}^K w_k \mathcal{R}^{(k)}(\mathcal{F}) + \sqrt{\sum_{k=1}^K \frac{(b-a)^2 w_k^2 \ln(1/\epsilon)}{2N_k}}.$$

This completes the proof. ■

### C.3 Proof of Theorem 6.2

In order to prove Theorem 6.2, we also need the following result (see Bousquet *et al.*, 2004, Theorem 5):

**Theorem C.3** *Let  $\mathcal{F} \subseteq [a, b]^{\mathbb{Z}}$ . For any  $\epsilon > 0$ , with probability at least  $1 - \epsilon$ , there holds that for any  $f \in \mathcal{F}$ ,*

$$\begin{aligned}
\mathbb{E} f &\leq \mathbb{E}_N f + 2 \mathcal{R}(\mathcal{F}) + \sqrt{\frac{(b-a) \ln(1/\epsilon)}{2N}} \\
&\leq \mathbb{E}_N f + 2 \mathcal{R}_N(\mathcal{F}) + 3 \sqrt{\frac{(b-a) \ln(2/\epsilon)}{2N}}. \tag{109}
\end{aligned}$$

Again, we prove Theorem 6.2 by using Theorems C.2 and C.3.

**Proof of Theorem 6.2** We only consider the result of Theorem C.2 in a special setting of  $K = 2$ ,  $N_1 = N_T$ ,  $N_2 = N_S$ ,  $w_1 = \tau$  and  $w_2 = 1 - \tau$ . Given a choice of  $\tau \in [0, 1)$ , denote

$$H(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) := \sup_{f \in \mathcal{F}} |\mathbb{E}_\tau f - \mathbb{E}^{(T)} f|. \quad (110)$$

By (9), we have

$$H(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) = \sup_{f \in \mathcal{F}} \left| \tau \mathbb{E}_{N_T}^{(T)} f + (1 - \tau) \mathbb{E}_{N_S}^{(S)} f - \mathbb{E}^{(T)} f \right|, \quad (111)$$

According to Theorem C.2, we have for any  $\xi > 0$ ,

$$\Pr \left\{ H(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) - \mathbb{E}^{(S)} \{ H(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) \} \geq \xi \right\} \leq \exp \left\{ \frac{-2\xi^2}{(b-a)^2 \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)} \right\},$$

which can be equivalently rewritten as with probability at least  $1 - (\epsilon/2)$ ,

$$\begin{aligned} & H(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) \\ & \leq \mathbb{E}^{(S)} \{ H(\mathbf{Z}_1^{N_S}, \bar{\mathbf{Z}}_1^{N_T}) \} + \sqrt{\frac{(b-a)^2 \ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)} \\ & = \mathbb{E}^{(S)} \left\{ \sup_{f \in \mathcal{F}} \left| \tau \mathbb{E}_{N_T}^{(T)} f + (1 - \tau) \mathbb{E}_{N_S}^{(S)} f - \mathbb{E}^{(T)} f \right| \right\} + \sqrt{\frac{(b-a)^2 \ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)} \\ & = \mathbb{E}^{(S)} \left\{ \sup_{f \in \mathcal{F}} \left| \tau (\mathbb{E}_{N_T}^{(T)} f - \mathbb{E}^{(T)} f) + (1 - \tau) (\mathbb{E}_{N_S}^{(S)} f - \mathbb{E}^{(T)} f) \right| \right\} \\ & \quad + \sqrt{\frac{(b-a)^2 \ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)} \\ & \leq \tau \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_T}^{(T)} f - \mathbb{E}^{(T)} f \right| + (1 - \tau) \mathbb{E}^{(S)} \left\{ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_S}^{(S)} f - \mathbb{E}^{(T)} f \right| \right\} \\ & \quad + \sqrt{\frac{(b-a)^2 \ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)} \\ & = \tau \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_T}^{(T)} f - \mathbb{E}^{(T)} f \right| + (1 - \tau) \mathbb{E}^{(S)} \left\{ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_S}^{(S)} f - \mathbb{E}^{(S)} f + \mathbb{E}^{(S)} f - \mathbb{E}^{(T)} f \right| \right\} \\ & \quad + \sqrt{\frac{(b-a)^2 \ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)} \\ & \leq \tau \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_T}^{(T)} f - \mathbb{E}^{(T)} f \right| + (1 - \tau) \mathbb{E}^{(S)} \left\{ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_S}^{(S)} f - \mathbb{E}^{(S)} f \right| \right\} \\ & \quad + (1 - \tau) \sup_{f \in \mathcal{F}} \left| \mathbb{E}^{(S)} f - \mathbb{E}^{(T)} f \right| + \sqrt{\frac{(b-a)^2 \ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)} \end{aligned}$$

$$\begin{aligned}
&= \tau \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_T}^{(T)} f - \mathbb{E}^{(T)} f \right| + (1 - \tau) \mathbb{E}^{(S)} \left\{ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_S}^{(S)} f - \mathbb{E}^{(S)} f \right| \right\} + (1 - \tau) D_{\mathcal{F}}(S, T) \\
&\quad + \sqrt{\frac{(b-a)^2 \ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)}. \tag{112}
\end{aligned}$$

According to Theorem C.3, for any  $\epsilon > 0$ , we have with at least  $1 - (\epsilon/2)$ ,

$$\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_T}^{(T)} f - \mathbb{E}^{(T)} f \right| \leq 2\mathcal{R}_{N_T}^{(T)}(\mathcal{F}) + 3\sqrt{\frac{(b-a) \ln(4/\epsilon)}{2N_T}}. \tag{113}$$

Next, according to (23), we have

$$\begin{aligned}
&\mathbb{E}^{(S)} \left\{ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_S}^{(S)} f - \mathbb{E}^{(S)} f \right| \right\} \\
&= \mathbb{E}^{(S)} \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_S}^{(S)} f - \mathbb{E}'^{(S)} \{ \mathbb{E}_{N_S}^{(S)} f \} \right| \\
&\leq \mathbb{E}^{(S)} \mathbb{E}'^{(S)} \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{N_k}^{(S)} f - \mathbb{E}'_{N_k}^{(S)} f \right| \\
&= \mathbb{E}^{(S)} \mathbb{E}'^{(S)} \sup_{f \in \mathcal{F}} \left| \frac{1}{N_S} \sum_{n=1}^{N_S} (f(\mathbf{z}_n^{(S)}) - f(\mathbf{z}'_n^{(S)})) \right| \\
&\leq \mathbb{E}^{(S)} \mathbb{E}'^{(S)} \mathbb{E}_{\sigma} \sup_{f \in \mathcal{F}} \left| \frac{1}{N_S} \sum_{n=1}^{N_S} \sigma_n (f(\mathbf{z}_n^{(S)}) - f(\mathbf{z}'_n^{(S)})) \right| \\
&\leq 2\mathbb{E}^{(S)} \mathbb{E}_{\sigma} \sup_{f \in \mathcal{F}} \left| \frac{1}{N_S} \sum_{n=1}^{N_S} \sigma_n f(\mathbf{z}_n^{(S)}) \right| \\
&= 2\mathcal{R}^{(S)}(\mathcal{F}). \tag{114}
\end{aligned}$$

By combining (110), (112), (113) and (114), we obtain with probability at least  $1 - \epsilon$ ,

$$\begin{aligned}
\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\tau} f - \mathbb{E}^{(T)} f \right| &\leq (1 - \tau) D_{\mathcal{F}}(S, T) + 2(1 - \tau) \mathcal{R}^{(S)}(\mathcal{F}) \\
&\quad + 2\tau \mathcal{R}_{N_T}^{(T)}(\mathcal{F}) + 3\tau \sqrt{\frac{(b-a) \ln(4/\epsilon)}{2N_T}} \\
&\quad + (1 - \tau) \sqrt{\frac{(b-a)^2 \ln(2/\epsilon)}{2} \left( \frac{\tau^2}{N_T} + \frac{(1-\tau)^2}{N_S} \right)}.
\end{aligned}$$

This completes the proof. ■